

Devanagari OCR using a recognition driven segmentation framework and stochastic language models

Suryaprakash Kompalli · Srirangaraj Setlur · Venu Govindaraju

Received: 11 January 2008 / Revised: 25 March 2009 / Accepted: 13 April 2009 / Published online: 19 May 2009
© Springer-Verlag 2009

Abstract This paper describes a novel recognition driven segmentation methodology for *Devanagari* Optical Character Recognition. Prior approaches have used sequential rules to segment characters followed by template matching for classification. Our method uses a graph representation to segment characters. This method allows us to segment horizontally or vertically overlapping characters as well as those connected along non-linear boundaries into finer primitive components. The components are then processed by a classifier and the classifier score is used to determine if the components need to be further segmented. Multiple hypotheses are obtained for each composite character by considering all possible combinations of the classifier results for the primitive components. Word recognition is performed by designing a stochastic finite state automaton (SFSA) that takes into account both classifier scores as well as character frequencies. A novel feature of our approach is that we use sub-character primitive components in the classification stage in order to reduce the number of classes whereas we use an n-gram language model based on the linguistic character units for word recognition.

1 Introduction

Optical Character Recognition (OCR), the process of automated reading of text from scanned images, is a vital tool for applications such as the creation of digital libraries. Latin and CJK (Chinese, Japanese, Korean) OCRs report high accuracies on machine-printed documents and have been integrated into several applications [12, 23, 32, 43]. The term OCR is

sometimes loosely used to include pre-processing steps such as binarization and noise removal, skew correction, and text block segmentation prior to actual recognition. The pre-processing methods employed by us have been dealt with in a previous article [17]. Several OCR challenges remain even after employing disciplined scanning methods and cutting-edge pre-processing techniques. Most *Devanagari* OCR methods segment a word into smaller units [2, 9, 30] for classification. Classifiers are designed to assign class labels to the segmented units. Post-processing using dictionary lookup and other techniques are used to improve recognition performance at the word level.

In this work, we focus on the recognition of printed *Devanagari* text after pre-processing, and propose a graph-based representation that facilitates a recognition-driven segmentation approach, and a stochastic finite state automaton for word recognition. Experimental results have been reported on the Indian Language Technologies (ILT) data set [20, 26]. The EMILLE Hindi text data corpus [1] has been used for generating the probabilistic n-gram language model.

Devanagari is a syllabic-alphabetic script [11] with a set of basic symbols—consonants, half-consonants, vowels, vowel modifiers, and special diacritic marks, which we will refer to as characters (Table 2). *Composite characters* are formed by combining one or more of these characters. We will refer to these as *composites*. The formation of composites from characters can be represented as a regular expression shown in Eq. (1); a similar form has been used previously in [2]. A horizontal header line (*Shirorekha*) runs across the top of the characters in a word, and the characters span three distinct zones (Fig. 1); an ascender zone above the *Shirorekha*, the core zone just below the *Shirorekha*, and a descender zone below the baseline of the core zone. Symbols written above or below the core will be referred to as ascender or descender components, respectively. A composite character

S. Kompalli · S. Setlur (✉) · V. Govindaraju
Department of Computer Science and Engineering,
University at Buffalo, State University of New York, Buffalo, USA
e-mail: setlur@buffalo.edu

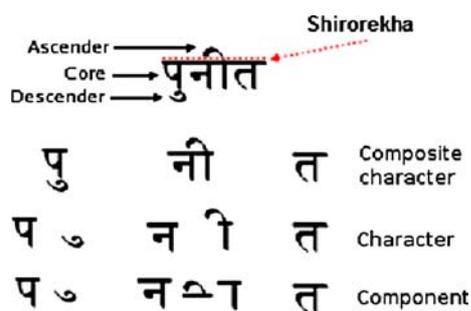


Fig. 1 Character zones in a *Devanagari* word (*top*) and nomenclature used for sub-elements of a word (*bottom*)

formed by one or more half consonants followed by a consonant and a vowel modifier will be referred to as a conjunct character or conjunct. Further, we use the following nomenclature for composite characters to indicate modifier location: Consonant–vowel (independent consonant or vowel with no ascenders or descenders), consonant–ascender (consonant with ascender), consonant–descender (consonant with descender), conjunct–core (conjunct with no ascenders or descenders), conjunct–ascender, and conjunct–descender.

$$\begin{aligned} \text{Composite character} &\rightarrow V|(c^*C[v][s]) \\ \text{where } V &\text{ - Vowel} \\ C &\text{ - Consonant} \\ c &\text{ - Half-consonant} \\ v &\text{ - Vowel modifier} \\ s &\text{ - Special accents} \end{aligned} \quad (1)$$

$$\text{Devanagari word} \rightarrow (\text{Composite character})^+$$

1.1 Prior work in segmentation

Removal of the *shirorekha* is a frequently used technique to split words into composites [9, 38]. Classification at the composite character level is very challenging as the number of valid classes exceeds 5,000. The ILT image data set [20, 26] has 10,606 word images containing 973 composite character classes, and the EMILLE text data set [1] has 5,573 composite character classes. Some sample images from the ILT data set are shown in Fig. 3. Typically, the difficulty of designing a classifier increases with the size of the class space, and finding sufficient samples of all classes for training becomes a critical problem. For example, 537 composite character classes from the ILT data set occur less than 25 times each. Since the number of such poorly represented classes is very large, OCR results are adversely affected.

An alternative approach [2, 17, 30] is to further segment the composite characters into smaller parts, which we will refer to as components. The class space of components is significantly smaller than that of composite characters. For example, the 11 composite characters in the set {ये यु या यि यो के कु का कि की को} are formed by

combining 8 characters: {य के ु ा ा ि ि ि ि}, and these in turn are formed from the following 5 components: {य के ु ा ि}. While there is some variation in the shape and location of components from one font to another, the basic shape and position of components remains fairly consistent. e.g.: The vowel modifier ि has an ascender part ु that always occurs above the shirorekha, and a part that occurs in the core region ा; However, minor variations in these components can be seen in the examples in Fig. 2(b).

The position of components may be estimated using structural attributes such as location and number of vertical strokes, average size, and aspect ratio. Structural information can be extracted from a paragraph or line of text, or based on knowledge of the font. However, font variations and poor print quality cause structural features to vary across words in the same line [24] (Fig. 2a, b). We show that it is possible to use a graph representation of a word image to capture the basic positions of components and also enable segmentation of composites. Structural information from paragraphs, lines, or a priori font data is not required to compute the graph representation.

Our work is partly inspired by the paradigm of recognition-driven OCR, originally proposed by Casey et al. [8]. *Devanagari* OCR methods reported in the literature [4] follow a dissection technique; a word is segmented into smaller parts, and classifiers are employed to assign class labels to the segmented image pieces. There is no feedback loop to correct erroneous segmentation. We design a stochastic framework that is used to prune multiple segmentation hypotheses.

Sliding window with forepart prediction has been extensively used in multi-font Latin OCR to separate fused characters [41]. A linear segmentation window is employed. Classifiers are used to map the two segmented parts to Latin characters. If the segmented parts produce classifier results with confidence less than a threshold, that segmentation hypothesis is rejected and the next hypothesis is considered. Linear windows sliding in only one direction do not work well with *Devanagari* script where composites need to be segmented along non-linear boundaries in two directions: left to right and top to bottom (Fig. 2c).

Lee et al. [29] outline a technique wherein they construct a graph by treating pixels of a gray-scale image as nodes that can be traversed to obtain segmentation hypotheses. While this method can trace a non-linear boundary, it leads to a large number of segmentation hypotheses.

1.2 Prior work in language modeling

Devanagari character formation can be said to be governed by a set of script composition rules [37]: e.g. Vowel modifiers should be associated with a consonant, only one vowel modifier must be present in a character. OCRs can

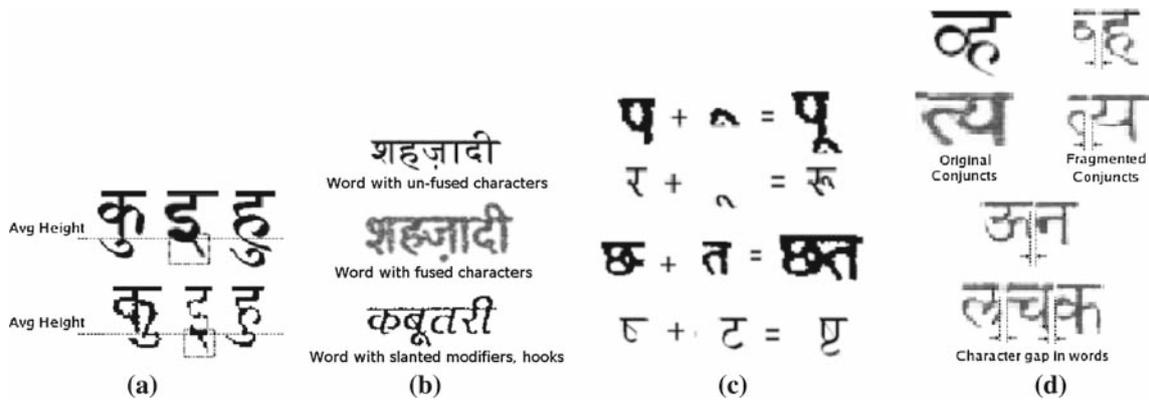


Fig. 2 Devanagari text samples: a descenders, b fused characters and stylized fonts, c composite characters, d fragmented conjuncts



Fig. 3 Sample images from the ILT data set

take advantage of these rules to improve performance [2, 30]. Dictionary lookup has been used as a post-processing step by Bansal et al. [3] to prune word recognition results. Use of probabilistic models which can enhance recognition results before using dictionary lookup have not been reported in the Devanagari OCR literature. Examples of such probabilistic models include alphabet n-grams which have been used in Latin OCR [7, 19], and word n-grams and trigger pairs which have been used in speech recognition systems [35]. Most of these probabilistic language models have used characters or words as units of the language model. At the sub-word level, attempts have been made to design models based on morphemes, or suffix-prefix relations [18, 39]. Readers should refer to Kukich [27] for an excellent survey of some of these techniques.

Sections 2.1 and 2.2 describe our recognition driven approach for segmentation and generation of multiple component hypotheses. Section 2.3 describes the design of a Stochastic Finite Automata (SFSA) that outputs word recognition results based on the component hypotheses and n-gram statistics. Section 2.4 describes the design of a language model based on the linguistic composite character units. Section 3 presents our experimental results with analysis. Section 4 outlines future work.

2 Proposed segmentation methodology

Let us define a graph representation G of a word or a character image I , derived from a topological representation [42] (Eq. 2).

$$G = (N, E)$$

where

$$N = \{N_i\} \text{ (Nodes, or Vertices)}$$

$$E = \{E_{ij}\} \text{ (Undirected edges)}$$

$$E_{ij} \begin{cases} 1 & \text{if there is an edge between } N_i, N_j \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

$$i, j = 1, 2, \dots, |N|$$

$$I = I_1 \cup I_2 \dots \cup I_{|N|}$$

$$I_i \cap I_j = \emptyset, \text{ if } i \neq j$$

Each graph G consists of a set of nodes N and a set of edges E . Each image I is divided into a number of horizontal runs (Hruns). Each of these runs are classified as *merging*, *splitting* or *continuing* runs based on the number of Hruns in the row immediately above it (*abv*) or in the row immediately below (*bel*) it. In *splitting* Hruns, $abv \leq 1$ and $bel > 1$ and in *merging* Hruns, $abv > 1$ and $bel \leq 1$. The remaining Hruns are classified as *continuing*. Adjacent runs which are classified as continuing are merged into a single block.

A splitting Hrun is added to the block that occurs above it and a merging Hrun is added to the block that occurs below it. A graph can be constructed by connecting the centroids of the neighboring blocks. A group of the *continuing* runs is defined as an image segment I_i or block. A node N_i of the graph G represents the center of mass of the image segment I_i . We define an edge E_{ij} when the corresponding nodes N_i, N_j are connected.

Variants of such a graph representation have been used in forms processing applications for form frame detection [46] and to extract features like loops, cusps, joints for classification in Chinese and Latin script recognition [34,45,47].

Figure 4 illustrates the process of generating a *block adjacency graph* (BAG) from a character image, and Fig. 5a shows the BAG generated from a word image.

After the BAG is constructed from a word image, baselines are computed using linear regression [40]. Blocks that lie parallel to the baseline and have a large aspect ratio and area greater than three times the average stroke width are identified as the shirorekha s . Blocks above the shirorekha corre-

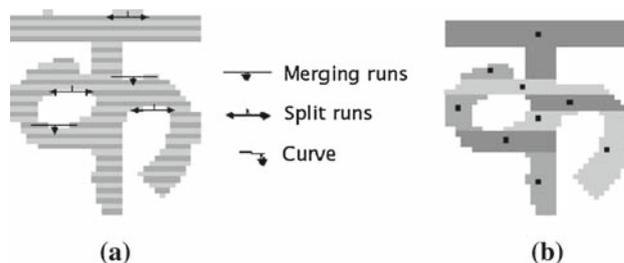


Fig. 4 Obtaining block adjacency graph. **a** Using horizontal runs (HRUNS) to identify initial blocks, **b** removing noisy HRUNS and creating blocks

spond to ascenders a , and the remaining blocks correspond to core components c . Since the shirorekha is very prominent, gross labeling of s, a, c is trivial even in noisy *Devanagari* words. The core components can be isolated characters that do not require further segmentation or conjunct and fused characters that may or may not have descenders. A recognition driven segmentation paradigm using graph features and classifiers is used to segment the conjuncts and fused characters.

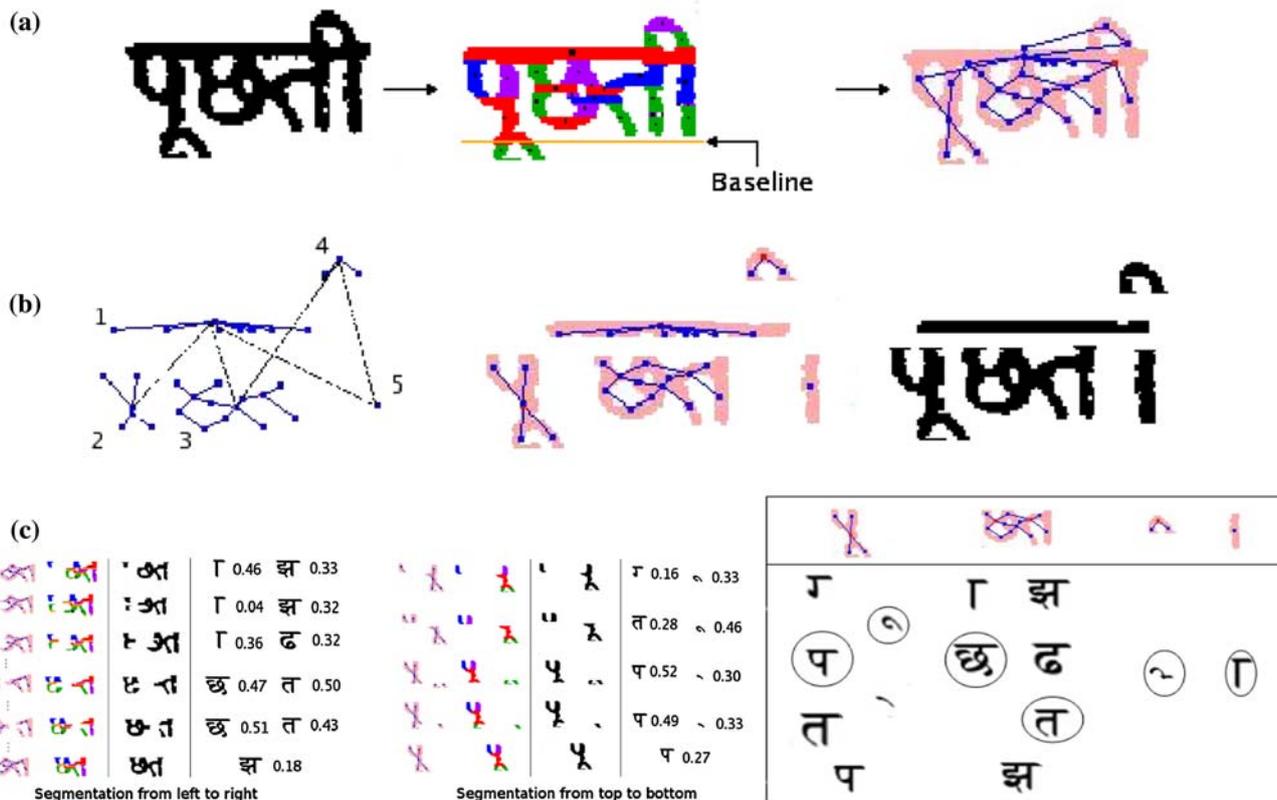


Fig. 5 Using block adjacency graphs (BAG) for Devanagari OCR. **a** Obtaining BAG from a word image; *left* original image, *center* BAG structure showing nodes, *right* edges between the nodes superimposed over the word image. **b** *Left* Edges selected for segmentation are shown in *dotted lines*; identification of shirorekha (1) and ascender (4), *center* preliminary segmentation hypothesis showing nodes for each segment,

right corresponding image segments. **c** Splitting a conjunct core component. *Columns 1, 2* Nodes of the graph selected as segmentation hypotheses, *column 3* image segments corresponding to the nodes, *column 4* classifier results for each image segment in column 3. **d** Lattice of classifier hypotheses. Results in circles are the correct class labels

2.1 Segmentation using graphs

We perform segmentation by selecting subgraphs from the BAG representation of word or character images. A sub-graph selected is a group of nodes from the BAG G and the ideal selection of sub-graphs will segment the image into its constituent *Devanagari* components. A naive approach is to exhaustively segment sub-graphs from the BAG. Such a scheme is computationally inefficient, since the number of possible hypotheses grows exponentially. The search space of hypotheses is pruned using the following constraints: (i) The primitives in the class space are assumed to be joined from left to right or top to bottom, and (ii) all of the primitives are assumed to consist of a single connected component. The equivalent conditions for sub-graph selection are: (i) A node cannot be selected if it is to the right (bottom for descender characters) of a node that has not been selected; and (ii) sub-graphs cannot have disconnected nodes. Using the graph representation, we can segment a composite character along curves, splits, and merges in the word image. The method generates fewer segmentation hypotheses than the window based approaches reported in literature [29,41]. Once the subgraphs have been extracted, they can be mapped back to the associated image segments thereby generating segmentation hypotheses. Images corresponding to the subgraphs are then input to a classifier (Sect. 2.2). Figure 5c shows an example of how different subgraphs can be selected from G to obtain segmentation hypotheses. Using the constraints described, the number of hypotheses can be further reduced.

2.2 Recognition driven segmentation

The use of classifiers to limit segmentation results has been used in Latin OCR [8,28]. We adopt a similar strategy for Devanagari. Table 1 shows the 71 component shapes that

Table 1 Class space of our system (components)

| | | | | | |
|---|---|---|-----|----|----|
| । | च | थ | र | क | ट |
| अ | छ | द | ल | रु | फ |
| इ | ज | ध | ळ | ह | ब |
| उ | झ | न | व | च | भ |
| ऊ | ट | प | श | ज | म |
| ए | ठ | फ | ष | झ | रु |
| क | ड | ब | स | त | ल |
| ख | ढ | भ | ह | थ | ठ |
| ग | ण | म | क्ष | ड | ढ |
| घ | त | य | त्र | न | रु |

Consonants and vowels Half-consonants

Ascenders Descenders

form our class space. There are 6 ascender shapes (a), 5 descender shapes (d) and 60 other core component shapes (c). For each subgraph extracted from the BAG of the word or character image, we input the corresponding subgraph image to a classifier. A score from the classifier below the threshold is taken as an indication that the input segmentation hypothesis does not correspond to one of the known classes. Procedure 1 outlines the specific steps used in our approach.

Procedure 1 Segment(G)

```

1: for all Word images  $w \in W$  do
2:   Start processing  $w$ 
3:   Obtain Block adjacency graph  $BAG(w)$  from  $w$ 
4:   Remove Shirerekha and detect baseline using  $BAG(w)$ 
5:   Generate multiple segmentation hypothesis  $H = \{h_1, h_2, h_3, \dots, h_k\}$  from  $BAG(w)$  via subgraph selection
6:   for all Segmentation hypothesis  $h_i \in H$  do
7:     if  $h_i$  represents an Ascender then
8:       Classify and remove ascender
9:     else
10:      Classify as core components
11:      if Classifier score > Threshold score obtained from ROC analysis then
12:        Call procedure RECOGNIZE( $h_i$ )
13:      else
14:        if Character crossed baseline then
15:          Segment character from top to bottom
16:        else
17:          Segment character from left to right
18:        end if
19:        Call procedure RECOGNIZE( $h_i$ )
20:      end if
21:    end if
22:  end for
23: end for

```

2.2.1 Classifiers

Classifiers are used to assign a class label from the class space for each input image and several classifier designs have emerged over the years. We have implemented previously known feature extraction techniques and applied them to Devanagari component classification. The features used in the experiments are gradient, structural and concavity features that are scale invariant and the classifiers are based on neural networks, or K-nearest neighbor. The classifiers can be considered as a function $score(I, C)$ where I is the input image, and C is the class space. The result S is of the format: $\langle o_{h_1}, s_{h_1} \rangle, \langle o_{h_2}, s_{h_2} \rangle \dots \langle o_{h_i}, s_{h_i} \rangle$, where o_{h_i} is the output class and $s_{h_i} \Rightarrow \mathbb{R}$ represents the class score or class-conditional probability obtained from the classifier.

1. Neural network classifier using Gradient features This classifier has 72 input nodes, 48 hidden nodes and 40 output nodes and uses the gradient sub-set of the GSC feature set. The GSC feature set was originally introduced for Latin handwriting recognition in [14]. The top-choice

Procedure 2 Recognize(h)

```

1: Initialize result set to return:  $R = \{\}$ 
2: Represent each hypothesis  $h$  as a tuple  $\langle h_1, h_2 \rangle$ 
3: Extract image  $I_{h_1}$  and  $I_{h_2}$  corresponding to hypothesis tuple
    $\langle h_1, h_2 \rangle$ 
4: for all  $\langle h_1, h_2 \rangle$  do
5:   if conjunct-core then
6:      $S_{h_1} = \text{score}(I_{h_1}, \text{valid core component})$ 
7:      $S_{h_2} = \text{score}(I_{h_2}, \text{valid core component})$ 
8:      $R = \{R, \langle S_{h_1}, S_{h_2} \rangle\}$ 
9:   end if
10:  if consonant-descender then
11:     $S_{h_1} = \text{score}(I_{h_1}, \text{valid core component})$ 
12:     $S_{h_2} = \text{score}(I_{h_2}, \text{valid descender component})$ 
13:     $R = \{R, \langle S_{h_1}, S_{h_2} \rangle\}$ 
14:  end if
15:  if conjunct-descender then
16:     $S_{h_2} = \text{score}(I_{h_2}, \text{valid descender component})$ 
17:     $R = \{R, \langle \text{conjunct-core} \rangle\}$ 
18:    Get recognition result for image  $I_{h_1}$  by calling
    RECOGNIZE( $h_1$ )
19:  for all  $\langle h_{11}, h_{12} \rangle$  do
20:     $R = \{R, \langle \langle h_{11}, h_{12} \rangle, S_{h_2} \rangle\}$ 
21:  end for
22: end if
23: end for

```

classifier accuracy on isolated consonant images is 94.36%

2. *Neural network classifier using GSC features* This classifier has 512 input nodes, 128 hidden nodes and 40 output nodes and uses the GSC feature set. The top-choice classifier accuracy on isolated consonant images is 95.25%.

3. *K-nearest neighbor classifier using GSC features* This is a K-nearest neighbor classifier that uses the GSC feature set and the Kulzinsky metric. This classifier has been previously used for handwriting recognition for Latin script [44]. We examined different values of K before selecting $K = 3$, which gave a top choice accuracy of 94.27% on isolated consonant images.

Both the neural networks are fully connected, and use logistic activation function. The neural networks have been trained using standard back-propagation. We retain the top 3 results returned by each classifier. The output scores for both the neural network as well as the K-nearest neighbor classifier range from 0.0 to 1.0. Receiver operator characteristics are analyzed and the equal error rate is selected as the threshold score for each class. Low thresholds increase the number of false positives (Conjuncts and descenders are accepted and under-segmentation occurs) and high thresholds increase false rejects (Consonants/vowels are rejected and over-segmentation occurs). The neural network classifier with GSC features provides the best result amongst the three methods.

The character recognition results from this stage are of the form: $O_h = \{o_{h_1} s_{h_1}, o_{h_2} s_{h_2}, \dots, o_{h_k} s_{h_k}\}$, where o and s represent the output class label and the associated confidence

score from the classifier. In cases where each hypothesis h_i represents a core component fused with another component (e.g. a conjunct character), the output is returned in the form of a tuple $\langle o_{h_i 1} s_{h_i 1}, o_{h_i 2} s_{h_i 2} \rangle$ where h_{i1} is the first part of the fused component and h_{i2} is the second part. In case of ascenders, only one hypothesis is obtained. For example, the results from the image shown in Fig. 5d can be written as follows:

$$\begin{aligned}
 O &= \{O_1, O_2, O_3, O_4\} \\
 O_1 &= \{\text{प}, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle\} \\
 O_2 &= \{\text{झ}, \langle \text{झ} \text{ र} \rangle, \langle \text{झ} \text{ त} \rangle, \langle \text{झ} \text{ ढ} \rangle, \langle \text{छ} \text{ र} \rangle, \langle \text{छ} \text{ त} \rangle, \\
 &\quad \langle \text{छ} \text{ ढ} \rangle\} \\
 O_3 &= \{\text{ँ}\} \\
 O_4 &= \{\text{ं}\}
 \end{aligned}$$

A probabilistic finite state automata is used to process these core, descender and ascender component hypotheses and generate word-level results.

2.3 Finite state automata to generate words

The recognition driven segmentation based on the BAG structure generates multiple hypotheses for each composite or fused character. While many words may be formed by combining these component results, not all of them are equally likely or valid. Script composition rules and corpus-based probabilistic language models can be used to eliminate invalid choices and rank the word results. We design a Stochastic Finite State Automata (SFSA) to combine the rule based script composition validity checking and a probabilistic n-gram language model into a single framework. We compare the use of various linguistic modeling units in selecting the n-grams.

The general idea in designing a SFSA is to create a Finite-State Machine (FSM) that accepts strings from a particular language, and assigning data-driven probabilities to transitions between states. If a candidate word is accepted, a score or probability is obtained and if rejected, corrections can be made using the minimum set of transitions that could not be traversed [21]. SFSA has been previously used to generate predictive text for hand-held devices [15], and to model error correcting parsers [22]. In both applications, the SFSA is built by studying n-grams using suffix/prefix relations in English text. SFSA has also been used in handwriting recognition to model the relation between image features like loops, cusps, or lines and the underlying alphabets and words [45]. In our methodology, the SFSA is modeled using script composition rules, and the transition probabilities are modeled on character or composite character n-grams. Some of the rules encode positional properties of components. For instance, descender components are placed below a core component (e.g. the modifier ँ is written below प to form पँ), or the vowel modifier ँ is made of the component ँ written on top of

Table 2 Decomposition of Devanagari characters into primitive components for recognition

Vowels having only one component:

अ इ उ ऊ ऋ ए

Vowels having multiple components:

आ → अ ा
 ई → इ ऀ
 ऐ → ए ँ
 ओ → आ ौ
 औ → ओ ौ | आ ौ

Consonants having only one component:

क ख घ च छ ज झ ञ ट ठ ड ढ त थ द ध न
 प फ ब भ म य र ल व ष स ह

Consonants having multiple components:

ग → ग ा
 ण → ण ा
 श → श ा

$$V \in \{अ आ इ ई उ ऊ ऋ ए ऐ ओ औ\}$$

(a) Vowel (V)

$$C \in \{क ख ग घ च छ ज झ ञ ट ठ ड ढ ण त थ द ध न प फ ब भ म य र ल व श ष स ह\}$$

(b) Consonant (C)

$$\text{half-consonant form, } f \in \{क ख ग घ च छ ज झ ञ ट ठ ड ढ ण त थ द ध न प फ ब भ म य र ल व श ष स\}$$

$$c \rightarrow f | C$$

(c) Half-consonant (c)

Modifiers having only one component:

ँ ॊ ो ौ ् ॎ ॏ ॐ ॑ ॒ ॓ ॔ ॕ ॖ ॗ क़ ख़ ग़ ज़ ड़ ढ़ फ़ य़ ॠ ॡ ॢ ॣ । ॥ ० १ २ ३ ४ ५ ६ ७ ८ ९ ॰ ॱ ॲ ॳ ॴ ॵ ॶ ॷ ॸ ॹ ॺ ॻ ॼ ॽ ॾ ॿ

Modifiers having two components:

ि → ि ि

ी → ि ि

े → ि ि

ै → ि ि

$$v \in \{ँ ॊ ो ौ ् ॎ ॏ ॐ ॑ ॒ ॓ ॔ ॕ ॖ ॗ क़ ख़ ग़ ज़ ड़ ढ़ फ़ य़ ॠ ॡ ॢ ॣ । ॥ ० १ २ ३ ४ ५ ६ ७ ८ ९ ॰ ॱ ॲ ॳ ॴ ॵ ॶ ॷ ॸ ॹ ॺ ॻ ॼ ॽ ॾ ॿ ि ि ि ि\}$$

(d) Vowel modifier (v)

Γ . The edges E of our graph representation (Eq. 2) inherently encode the connectivity between components. The recognition driven segmentation outlined in Procedure 1 translates this information into a sequence of class labels; For instance the class label sequence for पु and ि are: प्, ु, and ि respectively. Due to this inherent ordering present in class labels, positional information is not specified explicitly in our implementation of script composition rules.

Equation 1 represents a context sensitive grammar that describes the formation of valid Devanagari characters, and Table 2 describes the decomposition of characters into their sub-character primitives (components).

The SFSA is defined as follows: $\lambda = (Q, O, I, F, A, \Sigma, \delta)$, where

- Q : states
- O : observations (72 component classes, Table 1)

- I : initial states, which indicate start of a character
- F : final states, which indicate end of a character
- $A = a_{i,j}(o)$: set of transitions where $a_{i,j}(o)$ represents the probability of going from state i to j while accepting the label $o \in O$.
- Σ : emission symbols (characters from Table 2, and the null symbol ϵ)
- $\delta: A \times \Sigma$, a mapping from transitions to the emissions.

When a sequence of observations are accepted by the FSA, emissions generated by the corresponding transitions give the word result for this sequence. For example, the following transitions can accept the classifier labels given in Fig. 5d and

emit the correct word पूछती:

$$\begin{aligned}
 S &\xrightarrow{c} T_c \text{ emit } c, c \in \{\text{प छ त}\} \\
 T_c &\xrightarrow{d} T_d \text{ emit } d, d \in \{\epsilon\} \\
 T_c &\xrightarrow{m} NT_m \text{ emit } \epsilon, m \in \{\sim\} \\
 NT_m &\xrightarrow{l} T_v \text{ emit } v, l \in \{\Gamma\}, v \in \{\uparrow\} \\
 T_c, T_d, T_v &\xrightarrow{\epsilon} S \text{ emit } \epsilon
 \end{aligned}$$

where $S \in I$ and $T_c, T_d, T_v \in F$ and NT_m represents a non-terminal state. Emissions on the third and fourth transitions allow us to represent the combination of components \sim and Γ to form \uparrow . The sequence of emissions for this example would be: प छ ती. The four transitions in this example represent the following script writing grammar rules: (i) The modifier \uparrow is made of the ascender \sim connected to a consonant and a vertical bar, and (ii) Only one vowel modifier can be attached to a consonant. The resulting FSA is shown in Fig. 6.

In some SFSA models [45] the relation between states and observations is not apparent, and techniques like the Baum-Welch algorithm are used to learn transition probabilities. In our case, the observations are classifier symbols and the states are well-defined linguistic units indicating characters or composite characters. Using script composition rules, we parse the Devanagari text to obtain sequences of states for each word and the corresponding classifier symbols. Transition probabilities for the SFSA are obtained from text corpus using the formula:

$$a_{ij}(o) = \frac{\text{Count of observing } o \text{ going from state } i \text{ to } j}{\text{Number of transitions from state } i \text{ to } j} \quad (3)$$

Given a sequence of observation symbols, the best word hypothesis can be generated using the Viterbi algorithm. The Viterbi probability $\gamma_i(t)$, which is the highest probability of being in state i at time t is given by Eq. 4. The classifier score is also incorporated into this process (Eq. 5).

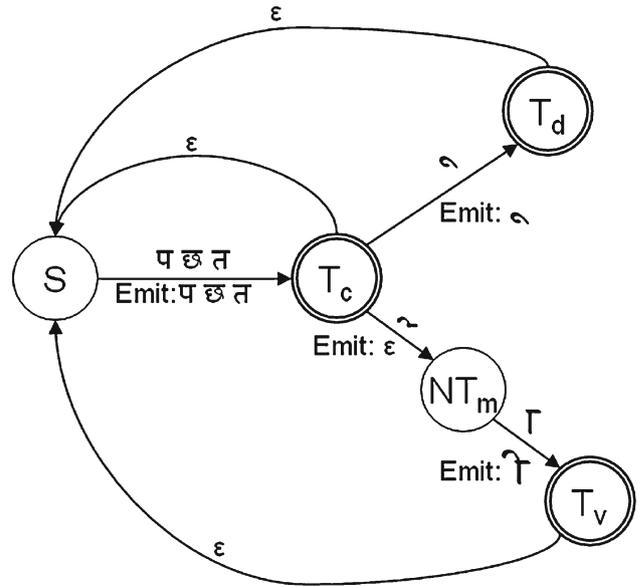


Fig. 6 Finite State Automata—the classifier labels in Fig. 5d are processed and the word पूछती: is emitted

$$\gamma_j(t) = \begin{cases} 1, & t = 0 \\ \max(\gamma_i(t-1)a_{ij}(o_t)), & \text{otherwise} \end{cases} \quad (4)$$

$$\gamma_j(t) = \begin{cases} 1, & t = 0 \\ \max(\gamma_i(t-1)a_{ij}(o_t)s_t), & \text{otherwise} \end{cases} \quad (5)$$

Tracing the transitions which generate $\gamma_j(T)$, we obtain the best word result corresponding to the matrix of component sequences (O_1, \dots, O_4) .

2.4 Hindi n-gram models

We now describe language specific inputs to the SFSA model to obtain multiple word hypotheses. To enhance the recognition results provided by the SFSA, we studied n-gram frequencies from the EMILLE Hindi corpus [1]. After removing punctuation, digits, and non-Devanagari letters the corpus has 11,381,720 words and 59,524,848 letters. One-third of the corpus was set aside for testing and the remaining two-thirds was used for training. An n-gram model specifies the conditional probability of each unit, or token of text c_n with respect to its history c_1, \dots, c_{n-1} [10]:

$$P(c_n|c_{1,\dots,n-1}) = \frac{P(c_1, \dots, c_n)}{P(c_1, \dots, c_{n-1})} = \frac{\text{Count}(c_1, \dots, c_n)}{\text{Count}(c_1, \dots, c_{n-1})} \quad (6)$$

We designed different n-gram models and report perplexity values over the test set (Pp):

$$Pp(D, q) = 2^{-\frac{1}{N} \sum_x \log q(x)} \quad (7)$$

where q represents the probability assigned by an n-gram model. A perplexity of k indicates that while predicting a

Table 3 Phonetic classification of Devanagari characters used for language modeling

| Voiceless Plosives | | Voiced Plosives | | Nasals |
|--------------------|-----------|-----------------|-----------|--------|
| Unaspirated | Aspirated | Unaspirated | Aspirated | |
| क | ख | ग | घ | ङ |
| च | छ | ज | झ | ञ |
| ट | ठ | ड | ढ | ण |
| त | थ | द | ध | न |
| प | फ | ब | भ | म |

token, the language model is as surprised as a uniform, independent selection from k tokens [10]. While language models with low perplexity are desirable, perplexity measures have to be examined in the context of the number of classes in the model. This is because perplexity decreases with decreasing number of classes. Consider a hypothetical model which has a single class “Devanagari character” and places all the tokens in our text into this category. Since this model places every token in the correct class (the class of “Devanagari character”), it will have a perplexity of zero. However, placing all tokens in a single group will not help an OCR. The token classes employed in an n-gram model are significant to model design and behavior. If two models have the same perplexity, it is likely that better modeling is being performed by the one having larger number of token classes. We experimented with models having different number of classes and studied the resulting variation in perplexity. Each model has a different binning approach where tokens of a specific type are included in a common class [10]. We have devised eight binning techniques for Devanagari and extract trigram and bigram frequencies. Some of our binning methods put characters having similar phonetic properties (Table 3) into a single bin. Motivation for such grouping comes from linguistic studies which state that consonants having similar phonetic properties are likely to exhibit similar co-occurrence behaviors [33]. In addition, each token c_i can be a character giving us letter n-grams, or tokens can be composite characters, giving us composite character n-grams.

1. *Characters with no classing* No classing is performed, i.e. each Devanagari character is considered as a separate class.
2. *Phonetic grouping of vowels and consonants* Tokens are the characters shown in Table 2 with the velar, palatal, retroflex, dental, labial consonants placed in five corresponding groups. Two groups are formed for sonorants, sibilants. Six groups each are formed for the vowels and vowel modifiers. Seven additional token classes have been used to account for special consonants (E.g. ज ख ग), digits, and punctuation marks.
3. *Aspirated plosive consonants in separate classes* Similar to case 2, except that each aspirated plosive consonant

Table 4 Bigram and trigram perplexity computed on the EMILLE data set

| Tokens | Number of classes | Perplexity | |
|---|-------------------|------------|---------|
| | | Bigram | Trigram |
| Characters | 126 | 13.64 | 6.23 |
| Characters grouped by phonetic properties | 26 | 8.28 | 6.06 |
| | 26 | 8.28 | 6.06 |
| | 36 | 9.00 | 6.10 |
| | 27 | 8.63 | 6.15 |
| | 28 | 9.00 | 6.23 |
| | 32 | 9.38 | 6.19 |
| | 41 | 9.71 | 6.15 |
| Composite characters | 5,538 | 18.51 | 3.53 |

(APC: ख घ छ झ ठ ढ थ ध फ भ) is placed in a group by itself.

4. *APCs in one class* Similar to case 2, except that all APCs are placed in a single class.
5. *Nasals in one class and APCs in separate classes* Similar to case 2, except all nasal consonants (NCs: ङ ञ ण न म) are grouped into a single class, and APCs in a different class.
6. *NCs in separate classes and APCs in a single class* Similar to case 2, except that each NC is assigned a separate class, and all APCs are placed in a single class.
7. *NCs and APCs in separate classes* Similar to case 2, except that each NC and APC is assigned to a separate class.
8. *Composite characters with no classing* Tokens correspond to the class of composite characters in the text.

We observe that with 5,538 classes, the composite character trigram and bigram perplexities are better than the character language model which has only 126 classes (Table 4). We have incorporated both character and composite character n-gram frequencies into the SFSA word generation model $\lambda = (Q, O, I, F, A, \Sigma, \delta)$.

2.4.1 Composite character n-grams in SFSA

The transition and Viterbi probabilities in Eqs. 3 and 4 do not capture character or composite character frequencies. While Eqs. 3 and 4 are defined with respect to each state Q of the SFSA, the terminal and start states respectively indicate which character or composite has been emitted. For example, while accepting classifier symbols for the word पूछती: (Fig. 6), the state sequences are: $S1 \rightarrow T_c2 \rightarrow T_d3 \rightarrow S4 \rightarrow T_c5 \rightarrow S6 \rightarrow T_c7 \rightarrow NT_m8 \rightarrow T_v9 \rightarrow S10$. In this example, $S1, S4, S6,$ and $S10$ are start states, $T_c2, T_d3, T_c5, T_c7, T_v9$ are terminal states, and NT_m8 is a non-terminal

state. To reach the terminal T_c2 , the FSA emits Ψ , indicating that the character Ψ is the recognition result. At the transition to T_d3 , it emits the character ∞ . The combination of these two characters gives the composite $\Psi\infty$. This composite is the emission between states $S1$ and $S4$. Similarly, ∞ and Ψ are emissions between $S4$ – $S6$ and $S6$ – $S10$. Transitions between the terminal and start states can capture character or composite n-gram statistics, but are not reflected in Eqs. 4 and 5. Therefore, we redefine the estimates as:

$$\hat{\gamma}_j(t) = \begin{cases} 1, & t = 0 \\ \max[\hat{\gamma}_i(t-1)a_{ij}(o_t)s_t], & \text{if } j \text{ is not a start state} \\ \max[\hat{\gamma}_i(t-1)a_{ij}(o_t)p(c_t|h_t)s_t], & \text{if } j \text{ is a start} \\ & \text{state, } c_t \text{ is the character between state } j \text{ and the} \\ & \text{previous start state, } h_t \text{ is the n-gram history of } c_t \end{cases} \quad (8)$$

To obtain the top- n results from our SFSA, we perform a recursive search through the lattice of component results and get all possible word recognition results ranked by $\hat{\gamma}_j(T)$ (Procedure 3). The time complexity of the search is exponential and would therefore be expensive for a generic graph. In practice, words have limited number of conjunct and descender hypotheses and an exhaustive search through the lattice of component results can be performed in real time. Dictionary lookup has been implemented using a lexicon of 24,391 words taken from the ILT data set. The lookup is performed using Levenshtein string edit distance with unit weights for substitution, deletion, and addition. The dictionary entry which gives minimum Levenshtien edit distance with SFSA output is taken as the dictionary corrected result.

Procedure 3 GetWords(O, λ)

1: Initialize: $t = 1$, Transitions vectors $V = V_1, \dots, V_n, \forall i$, set $V_i = \{ \} / * n = |O| * /$
 2: Initialize: Sequence of start states and emissions, $S = \{ \}$, Current emission $E = \{ \}$ and Position of current start state $pos = 0$
 3: Using Equation 8, find all $\hat{\gamma}_j(t)$ that accept O_1 , add them to V_1
 4: $t = t + 1$
 5: Parse(V, O, t, S, E, pos)

3 Experimental results and analysis

Our OCR was tested on 10,606 words from the University at Buffalo ILT data set. The data set contains images scanned at 300 dpi from varied sources such as magazines, newspapers, contemporary, and period articles. After binarization, line, and word separation, a test set of composite characters and components was extracted from word images using a semi-automated character segmentation technique [17]. The test set contains 8,209 samples of isolated consonants and vowels belonging to 40 classes, 686 samples of conjuncts belonging

Procedure 4 Parse($V, O, t, S, E, pos, \lambda$)

1: if V_{t-1} is not empty then
 2: if $t < N$ then
 3: Find $\hat{\gamma}_i(t-1) = \max[V_{t-1}]$
 4: Using equation 8, find all $\hat{\gamma}_j(t)$ that accept O_t , from state i , and add them to V_t
 5: $E = E +$ emission corresponding to $\hat{\gamma}_i(t-1)$
 6: if i is a start state then
 7: Add emission E to S
 8: Set $E = \{ \}$, $pos = pos + 1$
 9: end if
 10: Parse(V, O, t, S, E, Pos)
 11: else
 12: for all $\hat{\gamma}_i(t-1) \in V_{t-1}$ and $i \in F$ do
 13: word = $S + E +$ emission corresponding to $\hat{\gamma}_i(t-1)$
 14: score = $\hat{\gamma}_i(t-1)$
 15: Print result: (word, score)
 16: end for
 17: Backtrack($V, O, t-1, S, E, pos, \lambda$)
 18: end if
 19: else
 20: Backtrack($V, O, t, S, E, pos, \lambda$)
 21: end if

Procedure 5 Backtrack($V, O, t, S, E, pos, \lambda$)

1: if $t < 2$ then
 2: stop
 3: end if
 4: Find $\hat{\gamma}_i(t-2) = \max[V_{t-2}]$, remove $\hat{\gamma}_i(t-2)$ from V_{t-2}
 5: if i is a start state then
 6: $E =$ Last stored emission from S
 7: Delete last stored emission from S
 8: $pos = pos - 1$
 9: end if
 10: Parse($V, O, t-1, S, E, pos$)

to 97 classes, 479 samples of consonant–descenders belonging to 56 classes, and approximately 400 samples each of isolated ascenders and descenders. While many more valid conjunct and consonant–descender character classes can be obtained by combining the 38 consonants and 5 descenders, clean samples corresponding to only 153 (97+56) classes could be obtained from the data set.

3.1 Recognition of composite characters

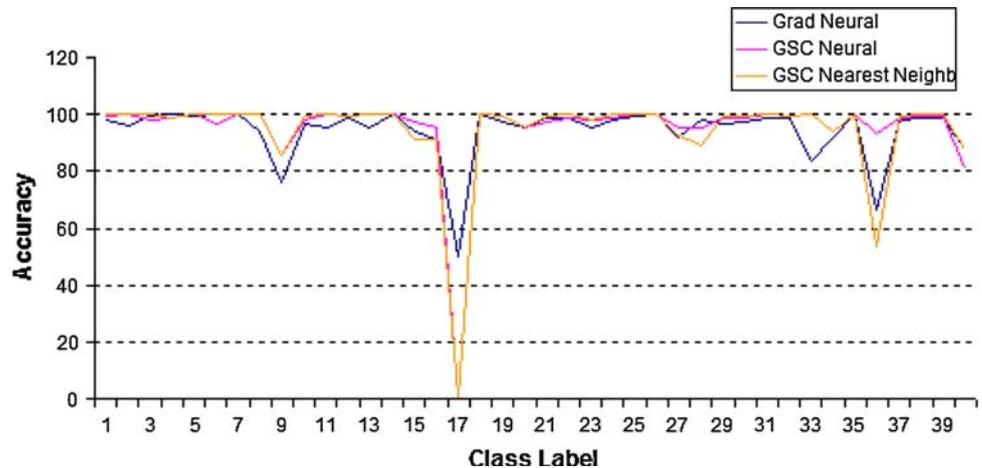
The accuracy of recognition of composite characters is analyzed with reference to the three procedural stages of our system: (i) accepting and recognizing an independent consonant/vowel and rejecting a conjunct/consonant–descender using the component classifier as indicated in Procedure 1, (ii) identifying the direction of segmentation for the rejected images as listed in Procedure 1, and (iii) using the graph to segment a composite from left to right or top to bottom and recognizing the constituent components as per Procedure 2.

In order to test stage (i), samples of consonants/vowels, conjunct and descender characters were input to the

Table 5 False reject rate (FRR) and false accept rate (FAR) obtained at the equal error rate for consonant/vowel classifiers, Testing is done on 8,029 consonants and vowels, 479 consonant–descenders and 686 conjuncts

| Classifier type | FRR (%) | FAR (%) | Breakup of FAR | |
|----------------------------------|---------|---------|-------------------------|--------------|
| | | | Consonant–descender (%) | Conjunct (%) |
| Neural network gradient features | 10.70 | 48.85 | 48.91 | 48.83 |
| Neural network GSC features | 5.83 | 6.88 | 7.40 | 6.70 |
| K-NN GSC features | 4.93 | 7.28 | 4.38 | 8.28 |

Fig. 7 Classifier accuracies for 8,029 consonant/vowel images belonging to 40 classes, Average accuracy: Neural network with gradient features: 94.36%, Neural network with GSC features: 95.25%, and K-nearest neighbor with GSC features: 96.97%



classifiers. False Accept Rate (FAR—a conjunct/descender is wrongly accepted), and False Reject Rates (FRR—a consonant/vowel is wrongly rejected) for different classifiers is reported in Table 5. The SFSA model is designed to use all classifier results and handle over-segmentation. False rejects can therefore be tolerated by the system better than false accepts. In our case, the best FAR and FRR rates are obtained using the GSC neural network classifier: 5.83% FRR and 6.88% FAR. Once a consonant/vowel sample is identified, it is classified into one of the consonant/vowel classes. The average accuracy of recognizing such consonants/vowels is in the range of 94–97% depending on the features used and the classification scheme (Fig. 7).

A majority of the composites that contribute to FAR are conjuncts that are formed by adding small modifiers to the consonant, for example, प्रफ्रप्र are formed by adding a small stroke (called rakaar) to the consonants प फ ग . The classifier confuses such rakaar-characters with the closest consonant. The most frequent rakaar character in our data set is प्र , and is frequently recognized as प with a high confidence. One solution to this problem is to include प्र in the class space of isolated core components and re-train the classifier. However, this still results in significant confusion between the two classes. This problem could also be potentially overcome during post-processing using a suitable language model.

In stage (ii), the system uses structural features from the graph to label the characters as either conjunct/consonant–

descenders or conjuncts. This stage has been tested using a data set containing 846 words which have at least one conjunct consonant, and 727 words having at least one conjunct–descender or consonant–descender. The operations are performed at the word level in order to estimate the baseline. Baselines are calculated for each word, and the characters are classified as conjuncts or conjunct/consonant–descenders. The typical errors in this step are due to the mislabeling of consonant–descenders like क्य as vertically joined conjunct characters.

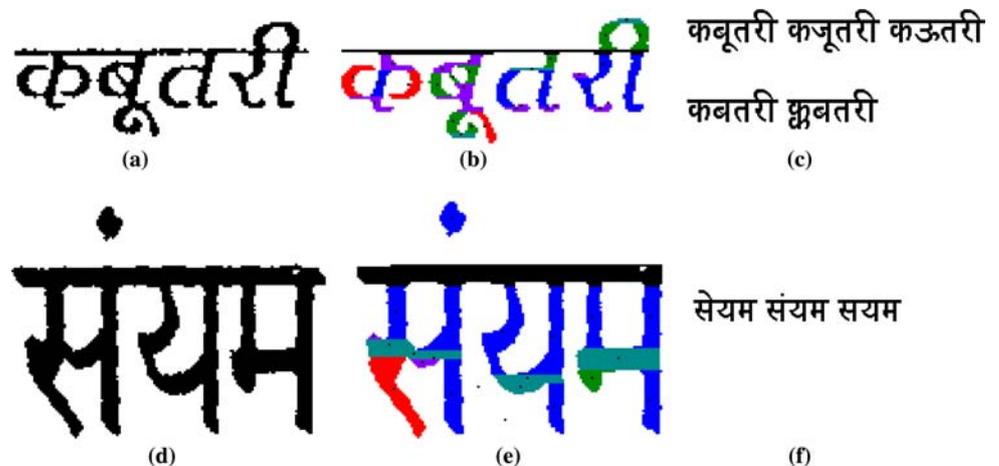
Results for stage (iii) are shown in Table 6. This stage has been tested using 686 conjunct character images and 479 conjunct/consonant–descender images. Top choice accuracy in recognizing these composites ranges from 26 to 38%, whereas the top-5 choice result ranges from 91 to 98%. Empirical analysis shows that errors are due to confusion between similar consonants/vowels; e.g. क्य being recognized as फ्य, फप , or कथ . Higher error rates are also seen in composites which contain three or more components, e.g. क्यु स्तू . However, such characters rarely occur in Devanagari text (0.47% in CEDAR-ILT, and 0.64% in EMILLE). Better modeling is required to overcome these issues at the character recognition level.

3.2 Word recognition

Word recognition performance was measured using four different stochastic finite state automaton (SFSA) models: (i)

Table 6 Character recognition results on 479 consonant–descender and 686 conjunct characters

| Character type | Recognized component | Top choice | | | | |
|--|----------------------|------------|-------|-------|-------|-------|
| | | 1 (%) | 2 (%) | 3 (%) | 4 (%) | 5 (%) |
| (a) K-Nearest neighbor classifier using GSC features | | | | | | |
| Consonant–descender | Consonant | 80.17 | 86.64 | 93.11 | 97.08 | 98.74 |
| | Descender | 84.76 | – | – | – | – |
| Conjuncts | Half-consonant | 16.91 | 51.46 | 67.49 | 83.97 | 94.17 |
| | Consonant | 46.50 | 59.91 | 71.43 | 82.80 | 92.27 |
| (b) Neural network classifier using GSC features | | | | | | |
| Consonant–descender | Consonant | 38.00 | 62.42 | 77.45 | 92.49 | 98.12 |
| | Descender | 86.31 | – | – | – | – |
| Conjuncts | Half-consonant | 25.66 | 52.48 | 74.20 | 82.94 | 85.57 |
| | Consonant | 27.41 | 47.10 | 66.62 | 83.10 | 91.11 |

Fig. 8 Results on words printed with different fonts; **a, d** input words, **b, e** graph representation, **c, f** word recognition result**Table 7** Word recognition accuracy and average string edit distance using different n-gram models. The test set has 10,606 words

| Top- <i>N</i> | Accuracy (%) and average string edit distance for SFSA models | | | | | | | |
|---------------------------------------|---|------|------------------|------|------------------|------|-------------------|------|
| | No n-gram | | Character bigram | | Composite bigram | | Composite trigram | |
| (a) Results without dictionary lookup | | | | | | | | |
| 1 | 13.75 | 1.83 | 51.12 | 1.09 | 69.28 | 0.80 | 64.76 | 1.06 |
| 3 | 50.77 | 1.07 | 66.99 | 0.73 | 75.19 | 0.61 | 76.68 | 0.70 |
| 5 | 63.32 | 0.77 | 71.63 | 0.63 | 76.79 | 0.55 | 77.77 | 0.65 |
| 20 | 75.99 | 0.49 | 76.50 | 0.49 | 77.88 | 0.50 | 78.05 | 0.58 |
| (b) Results with dictionary lookup | | | | | | | | |
| 1 | 34.44 | 2.38 | 61.03 | 1.48 | 74.96 | 1.03 | 72.66 | 1.20 |
| 3 | 68.42 | 1.19 | 77.19 | 0.96 | 82.64 | 0.78 | 82.23 | 0.89 |
| 5 | 80.29 | 0.81 | 81.78 | 0.78 | 85.04 | 0.68 | 83.72 | 0.80 |
| 20 | 88.37 | 0.50 | 87.74 | 0.54 | 87.74 | 0.55 | 85.33 | 0.72 |

Fig. 9 Word recognition accuracy on different documents; **a, b** proposed recognition driven segmentation approach, **c** segmentation driven approach [17]. The accuracy of segmentation driven approach varies widely across different documents, while the recognition driven segmentation approach provides more consistent accuracy

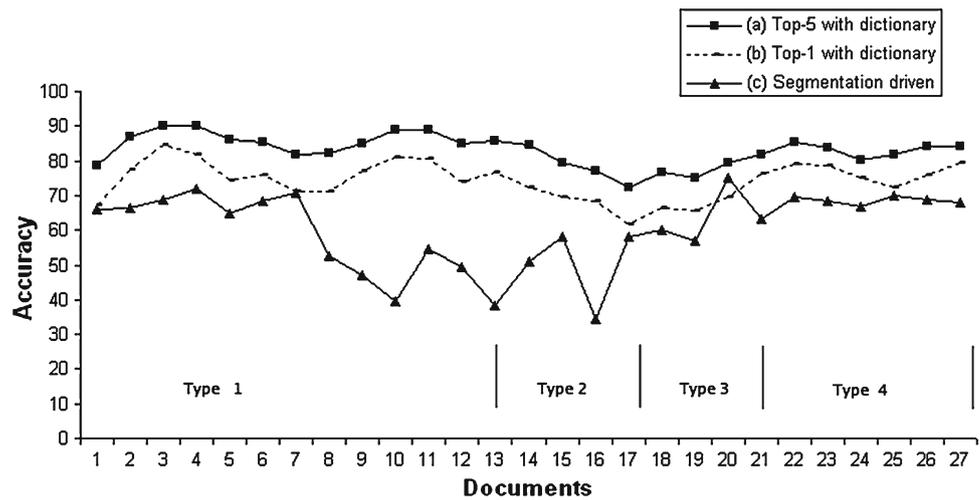


Table 8 Comparison of recognition results reported in the literature

| OCR Type | Classifier type | Data type and data set if known | Accuracy | |
|--------------------------------------|---------------------------|---|------------------------------|--------------|
| | | | Character (%) | Word |
| Segmentation driven OCR | Component classifier [2] | Two fonts | 93 | 67 |
| | Component classifier [30] | Single font (Oxford Hindi Dictionary) | 80 | 87 |
| | Character classifier [9] | Single font | 90 | 91 |
| | Component classifier [17] | CEDAR-ILT data set [26] | 84 | 53 |
| Recognition driven OCR | Current method, top-1 | CEDAR-ILT data set [26] | 96 | 75 |
| | Current method, top-20 | CEDAR-ILT data set [26] | 96 | 87 |
| Segmentation driven recognition [17] | | Graph based recognition driven segmentation | | |
| Image based recognition only | | Dictionary correction (%) | Image based recognition only | |
| Edit dist | (%) | | Edit dist | (%) |
| ≥ 4 | 29.11 +1.84 | | ≥ 4 | 10.9 +1.71 |
| 3 | 10.02 +2.24 | | 3 | 7.16 +2.72 |
| 2 | 10.52 +3.69 | | 2 | 12.67 +4.97 |
| 1 | 10.77 +5.53 | | 1 | 22.03 +18.64 |
| 0 | 39.58 | | 0 | 47.24 |
| Total: 2.36 | 100.00 | | 1.18 | 100.00 |
| Accuracy: | 52.89 | | | 75.28 |

Script writing grammar, with no n-gram statistics, (ii) SFSA with character bigrams, (iii) SFSA with composite character bigrams, and (iv) (SFSA) with composite character trigrams. A neural network classifier using GSC features was used for the experiments.

Table 6 summarizes the word recognition performance on the test set using the four models. Although the exhaustive search through the lattice (Procedures 3, 4, 5) is exponential in time complexity, it becomes tractable given the small number of conjunct and conjunct/consonant–descender hypotheses (less than ten). On average, one character generates 7

consonant–descender or conjunct hypotheses using our BAG approach. We achieved word level accuracy of 74.96% for the top choice using composite character bigram model with dictionary lookup. Examples of BAG generation and recognition on words of different fonts are shown in Fig. 8.

3.3 Comparison with prior methods

Our recognition driven segmentation paradigm is able to over-segment characters along non-linear boundaries when needed using the BAG representation. Classifier confidence

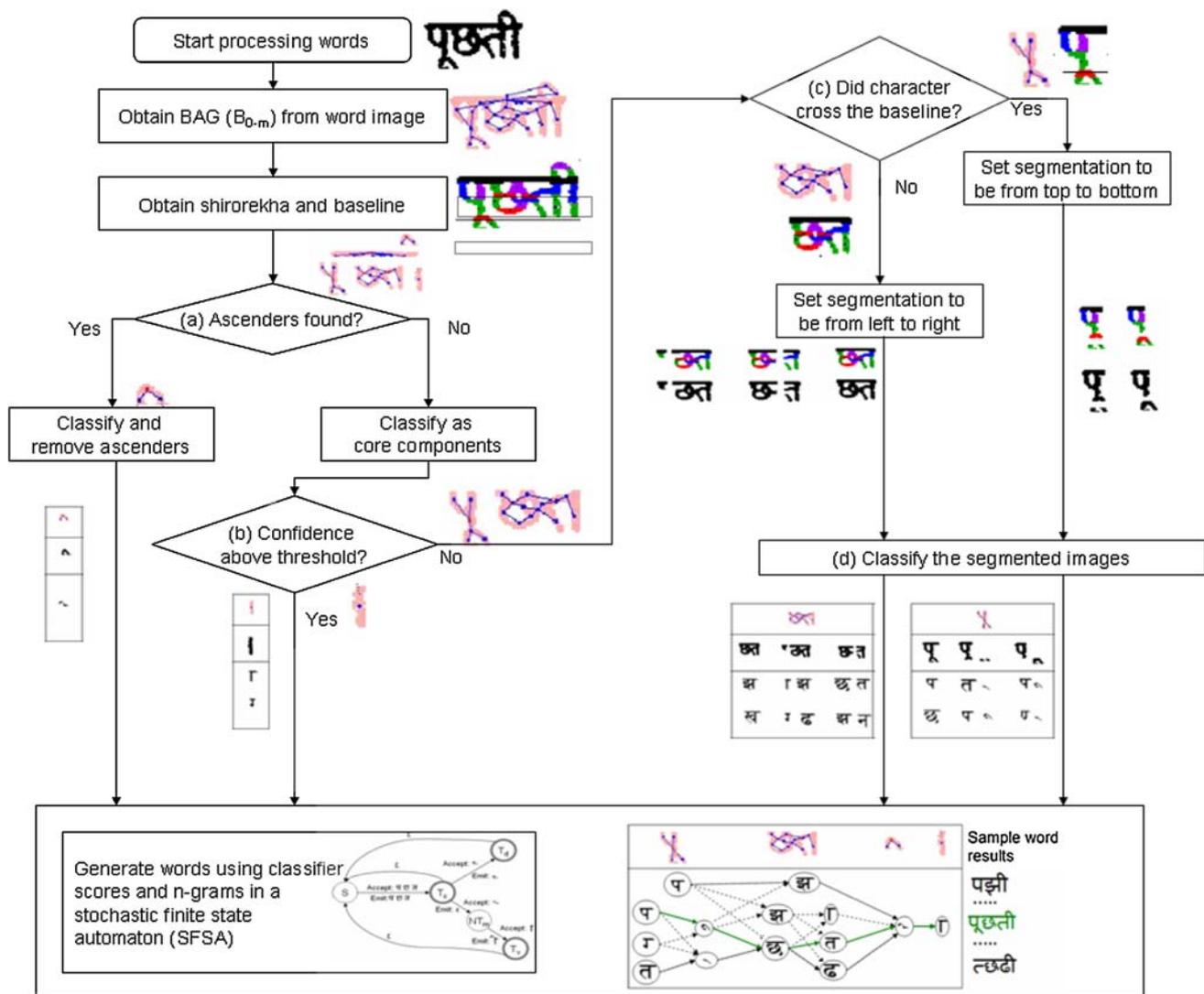


Fig. 10 Methodology overview

and language models are used to select the correct segmentation hypothesis. Previous techniques for conjunct character segmentation have used linear segmentation boundaries [5, 16]. Previous techniques for classification have used a picture language description (PLANG) of features [2, 36] and also decision trees [9]. These techniques use structural features specific to the *Devanagari* script such as vertical bars both during segmentation as well as subsequent classification. For comparison purposes, we used a recognizer developed in our lab using similar structural features [17, 24, 25].

The top-choice accuracy of the character classifier using structural features on the ILT test set is 74.53%, which is significantly higher than the top-choice accuracy of 26–38% obtained using the proposed method (Table 6). However, our method produces several segmentation hypotheses leading to a top-5 accuracy of 85–98%. We use linguistic information to reject the incorrect segmentation hypotheses. Given the

minor variations in character shapes between many classes, visual features alone are not sufficient. A language model is used to provide additional discriminating power. The results with and without dictionary lookup for the “Top-1, No n-gram” run are 32.99 and 13.25%, respectively (Table 7). Corresponding values for “Top-20, No n-gram” are 87.32 and 78.17%, and for “Top-1 with character bigram” the values are 74.71 and 70.17%. Figure 9 shows accuracy across different documents. The top-5 accuracy of the recognition driven BAG segmentation ranges from 72 to 90%, and the top-1 choice accuracy ranges from 62 to 85%. In comparison, top-1 accuracy of the segmentation driven approach ranges from 39 to 75%. Further, the segmentation driven approach does not provide multiple word results. Table 8 shows a comparison of results of Devanagari OCR systems reported in the literature. The recognition-driven OCR method shows significant gains over a purely segmentation-driven approach when measured

on the same data set. The character level accuracy of the current method is also higher than other results reported in the literature and comparable to results reported at the word level on other closed data sets. It has to be borne in mind that these numbers have been reported on different varied closed data sets.

4 Summary and future work

This paper presents a novel recognition driven segmentation methodology (Fig. 10) for Devanagari script OCR using the hypothesize and test paradigm. Composite characters are segmented along non-linear boundaries using the block adjacency graph representation, thus accommodating the natural breaks and joints in a character. The methodology used can readily accommodate commonly used feature extraction and classification techniques [6, 13, 31]. A stochastic model for word recognition has been presented. It combines classifier scores, script composition rules, and character n-gram statistics. Post-processing tools such as word n-grams or sentence-level grammar models are applied to prune the top- n choice results.

We have not considered special diacritic marks like *avagraha*, *udatta*, *anudatta*, special consonants such as ज्ञ ख़, punctuation and numerals. Symbols such as *anusvara*, *visarga* and the *reph* character often tend to be classified as noise. A Hindi corpus has been used to design the language model. Corpora from other languages that use the Devanagari script remain to be investigated.

Acknowledgments This material is based upon work supported by the National Science Foundation under Grants: IIS 0112059, IIS-0535038, IIS-0849511. We would like to thank Anurag Bhardwaj for helpful discussions and suggestions to improve the manuscript.

References

- Baker, P., Hardie, A., McEnery, T., Xiao, R., Bontcheva, K., Cunningham, H., Gaizauskas, R., Hamza, O., Maynard, D., Tablan, V., Ursu, C., Jayaram, B., Leisher, M.: Corpus linguistics and south asian languages: corpus creation and tool development. *Lit. Linguist. Comput.* **19**(4), 509–524 (2004)
- Bansal, V., Sinha, R.: Integrating knowledge sources in Devanagari text recognition. *IEEE Trans. Syst. Man Cybern. A* **30**(4), 500–505 (2000)
- Bansal, V., Sinha, R.: Partitioning and searching dictionary for correction of optically-read devanagari character strings. *Int. J. Doc. Anal. Recognit.* **4**(4), 269–280 (2002)
- Bansal, V., Sinha, R.: Segmentation of touching and fused Devanagari characters. *Pattern Recognit.* **35**, 875–893 (2002)
- Bansal, V., Sinha, R.: Segmentation of touching and fused Devanagari characters. *Pattern Recognit.* **35**, 875–893 (2002)
- Bishop, C.M.: *Neural Networks for Pattern Recognition*. Oxford University Press, New York (1996)
- Bouchaffra, D., Govindaraju, V., Srihari, S.N.: Postprocessing of recognized strings using nonstationary markovian models. *IEEE Trans. Pattern Anal. Mach. Intell.* **21**(10), 990–999 (1999)
- Casey, R., Lecolinet, E.: A survey of methods and strategies in character segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **18**, 690–706 (1996)
- Chaudhuri, B., Pal, U.: An OCR system to read two Indian language scripts: Bangla and Devanagari. In: *Proceedings of the 4th International Conference on Document Analysis and Recognition*, pp. 1011–1015 (1997)
- Christopher, M., Hinrich, S.: *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge (1999)
- Daniels, P.T., Bright, W.: *The World's Writing Systems*. Oxford University Press, New York (1996)
- Ding, X., Wen, D., Peng, L., Liu, C.: Document digitization technology and its application for digital library in china. In: *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries (DIAL 2004)*, pp. 46–53 (2004)
- Duda, R.O., Hart, P.E., Stork, D.G.: *Pattern Classification*, 2nd edn. Wiley, New York (2000)
- Favata, J., Srikantan, G.: A multiple feature/resolution approach to handprinted digit and character recognition. *Int. J. Imaging Syst. Technol.* **7**, 304–311 (1996)
- Forcada, M.: Corpus-based stochastic finite-state predictive text entry for reduced keyboards: application to catalan. In: *Procesamiento del Lenguaje Natural*, pp. 65–70 (2001)
- Garain, U., Chaudhuri, B.: Segmentation of touching characters in printed devnagari and bangla scripts using fuzzy multifactorial analysis. *IEEE Trans. Syst. Man. Cybern. C* **32**(4), 449–459 (2002)
- Govindaraju, V., Khedekar, S., Kompalli, S., Farooq, F., Setlur, S., Vemulapati, R.: Tools for enabling digital access to multilingual indic documents. In: *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries (DIAL 2004)*, pp. 122–133 (2004)
- Hirsimaki, T., Creutz, M., Siivola, V., Mikko, K.: Morphologically motivated language models in speech recognition. In: *Proceedings of International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning*, pp. 121–126 (2005)
- Hull, J.J., Srihari, S.N.: Experiments in text recognition with binary n-grams and viterbi algorithms. *IEEE Trans. Pattern Anal. Mach. Intell.* **4**(5), 520–530 (1982)
- The cedar-ilt data set. <http://www.cedar.buffalo.edu/ilt/>
- Juan, C.A., Enrique, V.: Efficient error-correcting viterbi parsing. *IEEE Trans. Pattern Anal. Mach. Intell.* **20**, 1109–1116 (1998)
- Juan, C.P.-C., Juan, C.A., Rafael, L.: Stochastic error-correcting parsing for OCR post-processing. In: *Proceedings of the 15th International Conference on Pattern Recognition*, vol. 4, pp. 405–408 (2000)
- Kim, G., Govindaraju, V., Srihari, S.N.: An architecture for handwritten text recognition systems. *IJDAR* **2**, 37–44 (1999)
- Kompalli, S., Nayak, S., Setlur, S., Govindaraju, V.: Challenges in ocr of devanagari documents. In: *Proceedings of the 8th International Conference on Document Analysis and Recognition*, pp. 327–333 (2005)
- Kompalli, S., Setlur, S., Govindaraju, V.: Design and comparison of segmentation driven and recognition driven Devanagari ocr. In: *Proceedings of the 2nd International Conference on Document Image Analysis for Libraries*, pp. 96–102 (2006)
- Kompalli, S., Setlur, S., Govindaraju, V., Vemulapati, R.: Creation of data resources and design of an evaluation test bed for Devanagari script recognition. In: *Proceedings of the 13th International Workshop on Research Issues on Data Engineering: Multi-lingual Information Management*, pp. 55–61 (2003)
- Kukich, K.: Techniques for automatically correcting words in text. *ACM Comput. Surv.* **24**(4), 377–439 (1992)

28. Kunihiro, F., Imagawa, T., Ashida, E.: Character recognition with selective attention. In: Proceedings of the International Joint Conference on Neural Networks, vol. 1, pp. 593–598 (1991)
29. Lee, S.-W., Lee, D.-J., Park, H.-S.: A new methodology for gray-scale character segmentation and recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* **18**, 1045–1050 (1996)
30. Ma, H., Doermann, D.: Adaptive hindi OCR using generalized hausdorff image comparison. *ACM Trans. Asian Lang. Inf. Process.* **26**(2), 198–213 (2003)
31. Mitchell, T.M.: *Machine Learning*. McGraw-Hill, New York (1997)
32. Mori, S., Suen, C.Y., Yamamoto, K.: Historical review of OCR research and development. *Proc. IEEE* **80**, 1029–1058 (1992)
33. Ohala, M.: *Aspects of Hindi Phonology*. Motilal Banarasidas, Delhi (1983). ISBN: 0895811162.
34. Rocha, J., Pavlidis, T.: Character recognition without segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **17**, 903–909 (1995)
35. Rosenfeld, R.: A maximum entropy approach to adaptive statistical language modeling. *Comput. Speech Lang.* **10**, 187–228 (1996)
36. Sinha, R.: Plang: a picture language schema for a class of pictures. *Pattern Recognit.* **16**(4), 373–383 (1983)
37. Sinha, R.: Rule based contextual post-processing for devanagari text recognition. *Pattern Recognit.* **20**, 475–485 (1987)
38. Sinha, R., Mahabala, H.: Machine recognition of Devnagari script. *IEEE Trans. Syst. Man Cybern.* **9**, 435–441 (1979)
39. Sinha, R., Prasada, B., Houle, G., Sabourin, M.: Hybrid contextual text recognition with string matching. *IEEE Trans. Pattern Anal. Mach. Intell.* **15**, 915–925 (1993)
40. Slavik, P., Govindaraju, V.: An overview of run-length encoding of handwritten word images. Technical report, SUNY, Buffalo (2000)
41. Song, J., Li, Z., Lyu, M., Cai, S.: Recognition of merged characters based on forepart prediction, necessity-sufficiency matching, and character-adaptive masking. *IEEE Trans. Syst. Man Cybern. B* **35**, 2–11 (2005)
42. Sonka, M., Hlavac, V., Boyle, R.: *Image Processing, Analysis and Machine Vision*, 2nd edn. Brooks-Cole, Belmont (1999)
43. Woo, K.J., George, T.R.: Automated labeling in document images. In: Proceedings of the SPIE, Document Recognition and Retrieval VIII, vol. **4307**, pp. 111–122, January 2001
44. Wu, Y., Ianakiev, K.G., Govindaraju, V.: Improved k-nearest neighbor classification. *Pattern Recognit.* **35**, 2311–2318 (2002)
45. Xue, H.: *Stochastic Modeling of High-Level Structures in Handwriting Recognition*. PhD thesis, University at Buffalo, The State University of New York (2002)
46. Yu, B., Jain, A.: A generic system for form dropout. *IEEE Trans. Pattern Anal. Mach. Intell.* **18**, 1127–1134 (1996)
47. Zheng, J., Ding, X., Wu, Y.: Recognizing on-line handwritten chinese character via farg matching. In: Proceedings of the 4th International Conference on Document Analysis and Recognition, vol. 2, pp. 621–624, August 1997