

ANALYSIS OF FORM IMAGES*

DACHENG WANG and SARGUR N. SRIHARI

CEDAR and Department of Computer Science

State University of New York at Buffalo

Buffalo, NY 14228-2567, USA

Email: srihari@cs.buffalo.edu

Automatic analysis of images of forms is a problem of both practical and theoretical interest; due to its importance in office automation, and due to the conceptual challenges posed for document image analysis, respectively. We describe an approach to the extraction of text, both typed and handwritten, from scanned and digitized images of filled-out forms. In decomposing a filled-out form into three basic components of boxes, line segments and the remainder (handwritten and typed characters, words, and logos), the method does *not* use a priori knowledge of form structure.

The input binary image is first segmented into small and large connected components. Complex boxes are decomposed into elementary regions using an approach based on key-point analysis. Handwritten and machine-printed text that touches or overlaps guide lines and boxes are separated by removing lines. Characters broken by line removal are rejoined using a character patching method. Experimental results with filled-out forms, from several different domains (insurance, banking, tax, retail and postal) are given.

Keywords: Segmentation, classification, connected component, key-point.

1. INTRODUCTION

A paper form is any preprinted document designed to elicit information. The problem of automatically segmenting an optically scanned and digitized image of a filled-out paper form containing typed and handwritten entries is considered in this study. The exact structure of the form is not assumed known a priori. The goal is to segment into the basic components of text (characters and words), line segments and boxes. The kinds of forms considered are those that contain preprinted boxes and the information elicited is entered using either handwriting or machine-printing (typing). In these forms, alphanumeric characters (machine-printed and handwritten) and certain special symbols play a key role in transmitting information between the individual filling out the forms and the individual requesting and reading the information. Forms contain preprinted queries and instructions eliciting information from the form-filler. The primary objective of segmentation is to extract all characters, preprinted and entered ones, so as to send them forward to a character recognition stage.

Several methods for segmenting a document into its components are described in the literature.¹⁻¹³ The methods aim to partition the document into text, graphics, half-tones, etc. These methods may be divided into divisive and agglomerative

*This work was supported in part by the National Science Foundation Grant IRI-86-13361, and by the Eastman Kodak Company.

methods (or top-down and bottom-up methods).¹⁴ Divisive methods are usually inappropriate for complex documents such as forms and engineering drawings. In such cases it is necessary to use an agglomerative method, one that determines connected components and groups them into larger units. The form analysis system proposed here involves a connected component analysis approach to segment characters, and a run length analysis approach to distinguish horizontal and vertical lines. The key-point method is used in this system to extract the structure of a form. We also propose a patching method to splice broken characters after removing touching line segments.

2. OVERVIEW OF FORM ANALYSIS SYSTEM

A form image analysis system needs to perform three preprocessing functions:

- (1) segment (or partition) into components and label them as boxes, line segments or text,
- (2) obtain spatial relationships between boxes and text,
- (3) separate text and lines which touch or overlap each other while avoiding breaking-up characters and words.

The data-flow in the form image analysis system is shown in Fig. 1. A binarized form image without skew is assumed to be prepared for input. Connected components are generated from the input image. These components are separated into small and large components. Small components contain individual characters, choice boxes, dotted lines, etc. Large components include not only line segments and complex boxes, but also touching character strings. For each large component, all black pixel run-lengths are separated into two groups: long run-length and short-run length. Each long run-length group is checked to determine whether it is a horizontal line segment candidate so as to leave it for further usage. A short run-length group may be a vertical line segment with or without touching characters, or it may be a character that touches other long run-length groups. So a step is needed to discriminate between vertical line segment candidates and short run-length groups. Another important step is to separate characters that touch line segments, and then splice these separated broken characters. A text map is obtained by merging small connected components, spliced characters as well as recognized character strings from large connected components.

The resulting horizontal and vertical line segment candidates are used to extract key-points that are form features. Twelve types of key-points are defined and associated with their coordinates: four end points, four convex corner points, and four concave corner points, which represent the end points and intersection points between horizontal and vertical line segments. Any area surrounded by several horizontal and vertical line segments can be described by a sequence of key-points. In other words, we can reconstruct a line segment frame from a list of key-points. The information of a form structure is obtained based on analyzing the relationship between key-points.

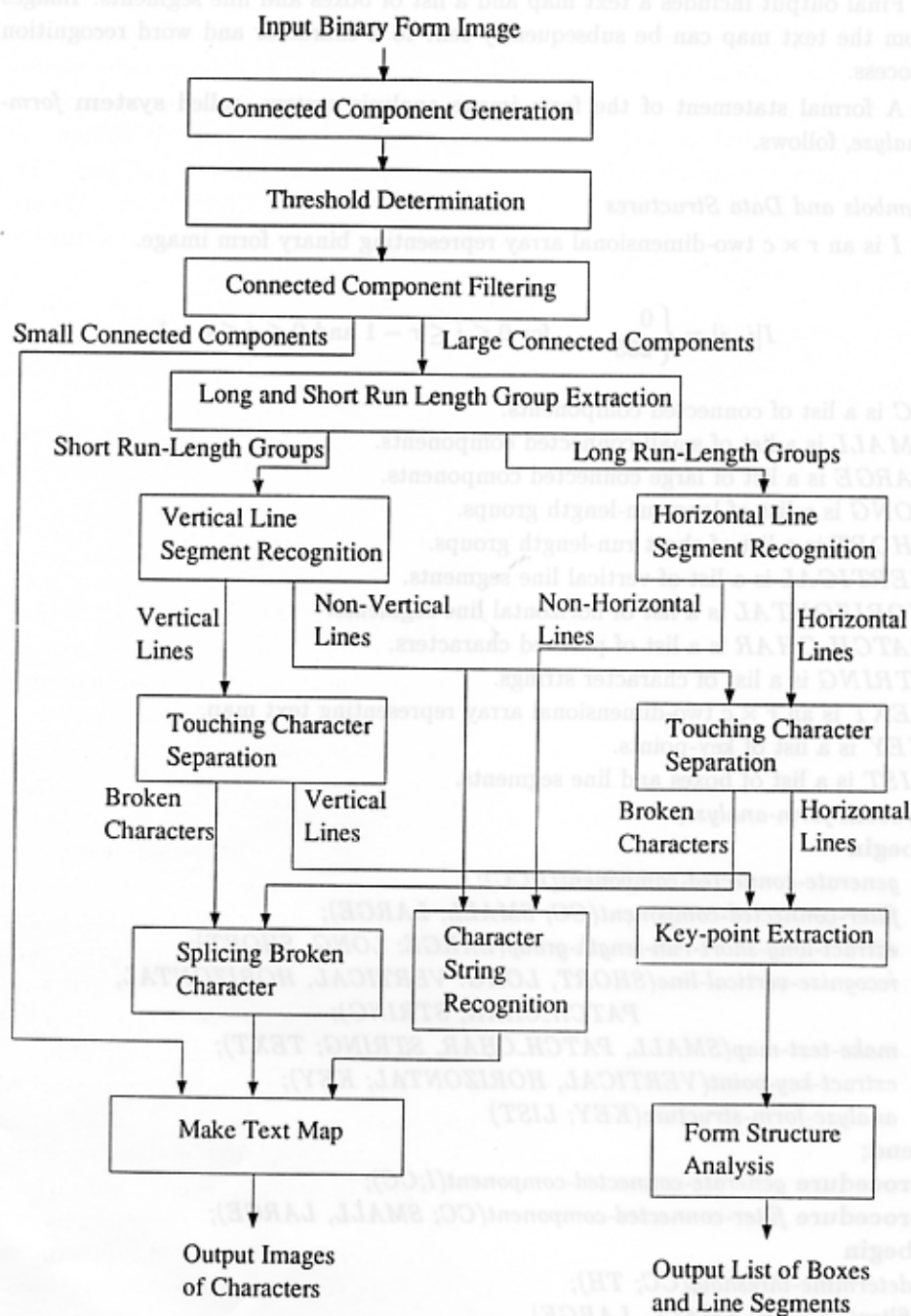


Fig. 1

Final output includes a text map and a list of boxes and line segments. Images from the text map can be subsequently sent to a character and word recognition process.

A formal statement of the form image analysis system, called *system form-analyze*, follows.

Symbols and Data Structures

I is an $r \times c$ two-dimensional array representing binary form image.

$$I[i, j] = \begin{cases} 0 \\ 255 \end{cases} \quad \text{for } 0 \leq i \leq r-1 \text{ and } 0 \leq j \leq c-1$$

CC is a list of connected components.

$SMALL$ is a list of small connected components.

$LARGE$ is a list of large connected components.

$LONG$ is a list of long run-length groups.

$SHORT$ is a list of short run-length groups.

$VERTICAL$ is a list of vertical line segments.

$HORIZONTAL$ is a list of horizontal line segments.

$PATCH.CHAR$ is a list of patched characters.

$STRING$ is a list of character strings.

$TEXT$ is an $r \times c$ two-dimensional array representing text map.

KEY is a list of key-points.

$LIST$ is a list of boxes and line segments.

system form-analyze;

begin

generate-connected-component($I; CC$);

filter-connected-component($CC; SMALL, LARGE$);

extract-long-short-run-length-group($LARGE: LONG, SHORT$);

recognize-vertical-line($SHORT, LONG: VERTICAL, HORIZONTAL,$
 $PATCH.CHAR, STRING$);

make-text-map($SMALL, PATCH.CHAR, STRING; TEXT$);

extract-key-point($VERTICAL, HORIZONTAL; KEY$);

analyze-form-structure($KEY; LIST$)

end;

procedure *generate-connected-component*($I; CC$);

procedure *filter-connected-component*($CC; SMALL, LARGE$);

begin

determine-threshold($CC; TH$);

filter($CC, TH; SMALL, LARGE$)

end;

procedure *extract-long-short-run-length-group*($LARGE; LONG, SHORT$);

begin

repeat

```

    group-long-run-length(LARGE; LONG);
    group-short-run-length(LARGE; SHORT)
until end-of-file
end;
procedure recognize-vertical-horizontal-line(LONG, SHORT; VERTICAL,
      HORIZONTAL, PATCH_CHAR, STRING);
begin
    recognize-vertical-line(SHORT; VERTICAL, PATCH_CHAR);
    recognize-horizontal-line(LONG; HORIZONTAL, PATCH_CHAR, STRING)
end;
procedure make-text-map(SMALL, PATCH_CHAR, STRING; TEXT);
procedure extract-key-point(VERTICAL, HORIZONTAL; KEY);
procedure analyze-form-structure(KEY; LIST)

```

3. SEGMENTATION INTO TWO PARTITIONS

The initial steps in form analysis are to determine connected components and to separate them into small and large components.

Connected component generation

For the purpose of defining digital connectivity, we assume 8-connectivity for 1-pixels (black pixels) and 4-connectivity for 0-pixels (white pixels). All 8-connected 1-pixels belonging to individual characters or boxes are grouped together and enclosed in bounding rectangles. Each rectangle thus represents an area in which a single connected component exists. A large box is enclosed by a large rectangle which may overlap other small rectangles. An enclosing rectangle corresponding to a character usually contains a complete character except in the case of characters such as a lowercase "i" which needs two enclosing rectangles. A connected component may also contain more than one character if characters touch their neighbors.

The algorithm for generating connected components is based on the line adjacency graph method of Pavlidis.¹⁵ However, we create a line adjacency data structure instead of a graph. This data structure records all black pixel run-lengths as well as their linkages to the adjacent runs located above and below. A connected component is extracted by tracing a group of runs that are linked together.

Connected component filtering

Both the height and width of a connected component corresponding to a character are generally smaller than those of a box. Either the height or width of a character component is smaller than that of a line segment. The objective of connected component filtering is to perform a preliminary separation of characters and boxes as well as line segments. Both the height and width of each component are compared with a threshold that is slightly larger than the largest point size of characters appearing in the form image. A component is classified as a character if both height and width are lower than the threshold. Certain other small connected

components, e.g., boxes to be checked, dotted lines, noise etc., are also classified as characters since size is the only feature used.

The filtering threshold is automatically determined, as follows. A histogram of connected component heights is determined. Because the number of characters is much larger than the number of boxes and line segments, a large peak corresponding to characters appears in the histogram. The filtering threshold is estimated as the size where 95 percent of the cumulative distribution of the size histogram is contained.

4. DECOMPOSITION OF COMPLEX BOXES

A box is composed of several horizontal and vertical line segments. If we know the location and type of the end and intersection points that exist in all line segments, then we can obtain the structure information of the box. This is the basic idea of the key-point method. In this section we define key-points and describe how key-points are used to express boxes.

Key-point definition

A line segment is considered as having some thickness. That means any line segment in a form image consists of a group of black pixels which has a rectangular shape. Generally the height of a horizontal line segment or the width of a vertical line segment is very small. For simplicity, in defining a key-point, we ignore two short edges of the line segment so a line segment can be considered as a group of black pixels associated with two long edge lines beside its two sides as shown in Fig. 2(a) and (b). The end of a long edge line is referred to as an end key-point. Four types of end key-points are defined according to their direction: left-end, right-end, up-end, and down-end. A horizontal line segment has two left-end key-points and two right-end key-points, whereas a vertical line segment has two up-end key-points and two down-end key-points.

When several horizontal and vertical line segments touch or overlap together, the intersection points of two edge lines form another kind of key-point, viz., a corner key-point. Eight types of corner key-points are defined according to their position. Four are convex corners, and four are concave corners. Figures 2(c) and (d) illustrate the concave corner and convex corner key-points respectively. Table 1 gives definitions of all twelve key-points. Each key-point has x and y coordinates associated with it so as to point out its position in the form image.

Key-point expression of boxes

The structure of a box composed of horizontal and vertical line segments can be expressed by several sequences of key-points each of which corresponds to an individual region of the box.

A complex box consists of several regions which do not share a common area with each other. Two region types are defined: closed and open. A closed region is formed consisting of several horizontal and vertical line segments closely surrounding this

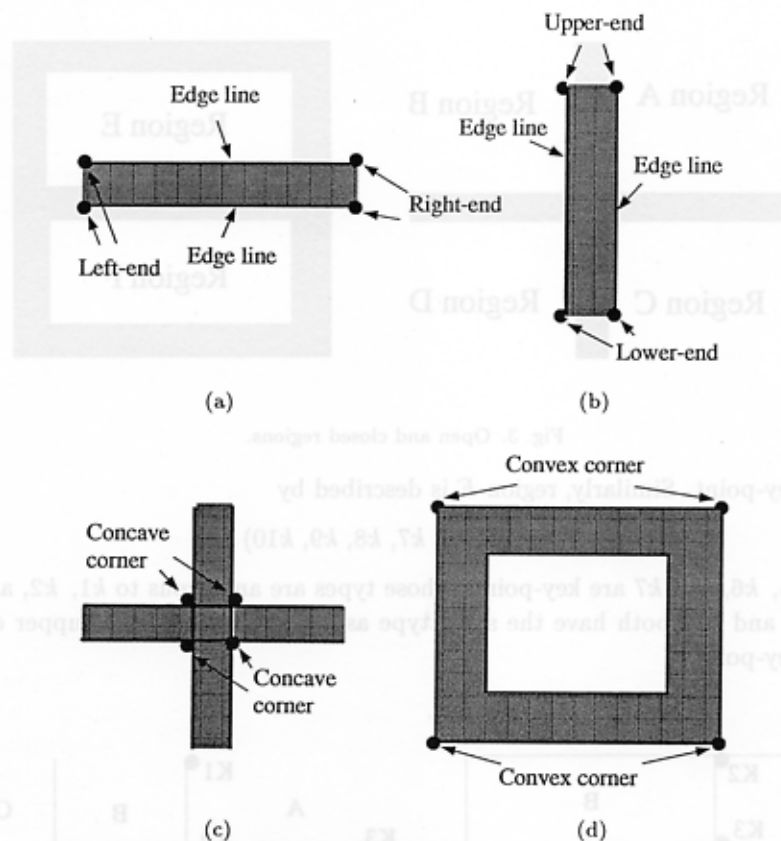


Fig. 2. Key-point definition. (a) Horizontal line segment and end key-point. (b) Vertical line segment and end key-point. (c) Concave corner key-point. (d) Convex corner key-point.

region. An open region has several horizontal and vertical line segments surrounding it, but these line segments do not form a closed boundary line. Examples of closed and open regions are given in Fig. 3, where regions *E* and *F* are closed, and regions *A*, *B*, *C*, and *D* are open.

The shape and position of a region can be determined by a sequence of key-points. For closed regions, there is no end key-point in the expression sequence. End key-points do appear in the expression sequence of key-points for an open region. The description of a region by a sequence of key-points is illustrated with the example box shown in Fig. 4: the complex box with closed boundary in (a) has seven closed regions, *A* to *G*. Region *A*, for example, is described by the following sequence of key-points:

$$A = (k1, k2, k3, k4),$$

where *k1* is a left upper concave corner key-point, *k2* is a right upper concave corner key-point, *k3* is a right lower concave corner key-point, and *k4* is a left lower concave

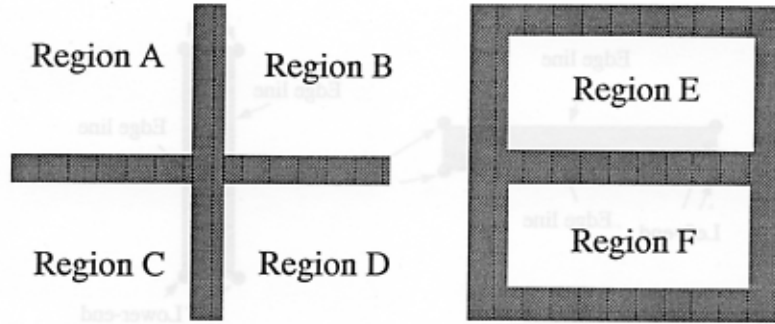


Fig. 3. Open and closed regions.

corner key-point. Similarly, region F is described by

$$F = (k5, k6, k7, k8, k9, k10),$$

where $k5$, $k6$, and $k7$ are key-points whose types are analogous to $k1$, $k2$, and $k3$, while $k8$ and $k10$ both have the same type as $k4$, and $k9$ is a right upper convex corner key-point.

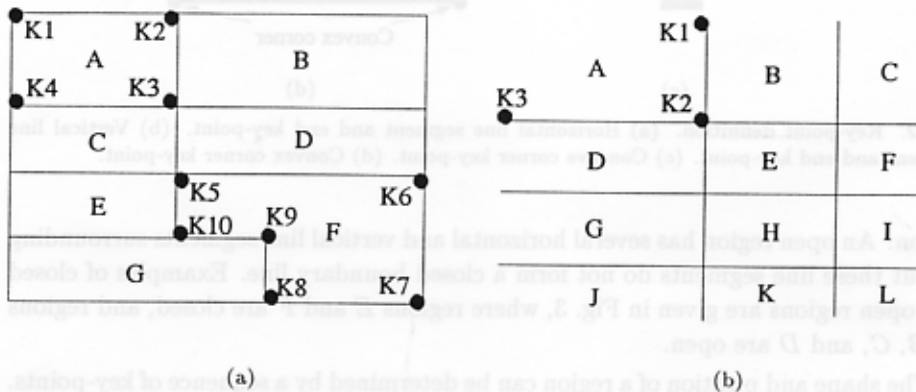


Fig. 4. Neighboring key-points: (a) regions included in a box with closed boundary: $k1$ is the horizontal neighbor of $k2$, and $k3$ is the vertical neighbor of $k2$, (b) regions included in a box without closed boundary: $k2$ is the vertical neighbor of $k1$, and $k1$ does not have a horizontal neighbor because it is an up end key-point.

A complex box that does not have a closed boundary is shown in Fig. 3(b). It consists of only open regions such as A , B , C , D , F , G , I , J , K , and L . Region A , for example, can be described by

$$A = (k1, k2, k3),$$

where $k1$ and $k3$ are end key-points, and $k2$ is a concave corner key-point.

Key-point extraction

In order to extract key-points from a form image, we have to first determine horizontal and vertical line segments. The horizontal and vertical line segments can be found from the line adjacency data structure obtained in the connected component generation stage.

For all large connected components, their run-lengths are checked again. Adjacent long run-lengths, which are longer than the filtering threshold used in the connected component filtering stage, are grouped together to become a candidate for a horizontal line segment. This approach determines real horizontal line segments easily. However, it also determines some "pseudo" horizontal line segments in character string areas. So a step to remove these pseudo-horizontal lines is needed, as described later in this section.

Adjacent short run-lengths, which are shorter than the filtering threshold used in the connected component filtering stage, are grouped together to form a short run-length region. A short run-length region may be (i) a character (or character string) that touches a candidate for a horizontal line segment, (ii) a vertical line that may touch a candidate for a horizontal line segment, or (iii) a vertical line with some touching characters. The projection profile is used to discriminate vertical line segments from short run-length regions. The projection profile in the vertical direction is determined for each short run-length region. The peak in the resulting projection profile is checked using two rules:

- (1) the value of the peak is equal to the height of the short run-length region, and
- (2) the width of the peak region is less than the filtering threshold.

A vertical line segment is found if the peak rules are satisfied for the projection profile of a short run-length region. The width of the vertical line segment is equal to the width of the peak region.

After obtaining vertical line segments, all candidate horizontal line segments are checked again to remove pseudo-horizontal lines found in character string areas. For these pseudo-horizontal lines, there are some short run-length regions touching the two ends of the horizontal line, but generally these short run-length regions do not have vertical line segments in them. This fact suggests setting up a rule for removing pseudo-horizontal lines.

Key-points are extracted easily by checking intersection points and end points of horizontal and vertical line segments. The end of a horizontal or vertical line segment forms a pair of end key-points if the end does not touch anything else. The intersection point of two line segments forms few corner key-points. The linkage between adjacent key-points is also extracted as the key-point is extracted. For each key-point, two neighbors — one each in the horizontal and vertical directions — are searched. If two key-points lie on the same horizontal or vertical line segment, then they are horizontal or vertical neighbors to each other, respectively (see Fig. 4).

Searching regions

Regions are searched after obtaining the list of key-points. Any key-point can be chosen as the starting point for the search. If this key-point has a horizontal neighbor, then the initial search direction is set as horizontal. Otherwise the initial search direction is set to be vertical. The search will move to the neighbor key-point in the initial search direction. Then the new search direction, different from the current search direction, is determined. Moving to the neighboring key-point and changing searching direction will be executed recursively. There are two cases that terminate search. If search reaches the start point again, then search is terminated, and a closed region is found. If search reaches an end key-point, then the search is paused at this end point and search is restarted from the start point again in an alternative direction. Search will terminate at another end key-point, and an open region is found.

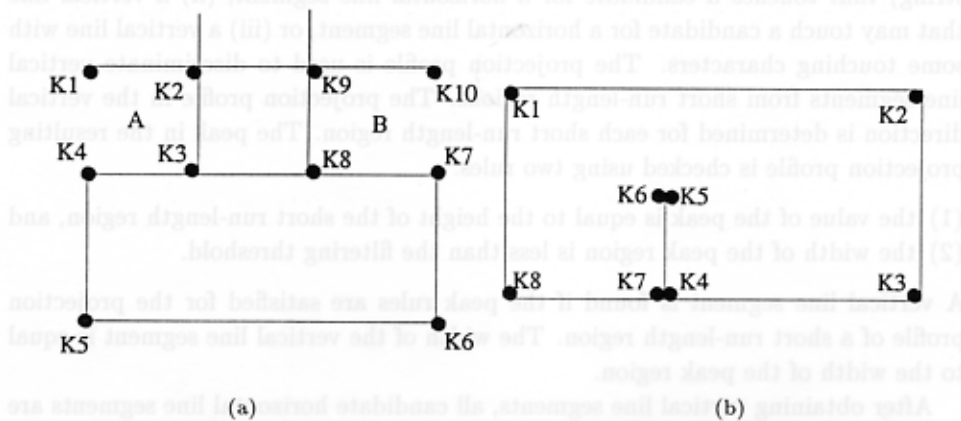


Fig. 5. Two cases of box structure for which two additional rules are created to describe regions.

Two exceptions need to be handled (see Fig. 5). The first is shown in (a) if we start search from end key-point k_1 , then search will be terminated at another end key-point k_{10} . However, the sequence of key-points obtained includes two open regions A and B as well as part of the outer border of the lower box. So an additional rule for this case is that if search starts from an end key-point and hits two sequential convex corner key-points, then terminate search at the first convex corner key-points. Region A is described as (k_1, k_2, k_3, k_4) , and region B as (k_{10}, k_9, k_8, k_7) . Another sequence of key-points (k_5, k_6) just describes part of the outer border of the lower rectangular box.

The second exception is shown in (b), where a few vertical line segments touch a closed box. The sequence of key-points for this region is $(k_1, k_2, k_3, k_4, k_5, k_6, k_7, k_8)$ where k_5 and k_6 are up end key-points. We decide that this region is closed although it has two (or more) up end key-points in its key-point expression sequence.

5. SEPARATING CHARACTERS CONNECTED TO LINES

Text that is in physical contact with a box or line segment remains in large connected components after connected component filtering. It is a necessary segmentation step to separate each character from the box or line segment it touches. Two cases of character-line contact are distinguished (see Fig. 6) (a) *touching*: the end of a character stroke touches a line segment, and (b) *overlapping*: character strokes cross over a line segment. In the touching case, the strokes of a given character have only one pair of intersection points with the line segment. In the overlapping case the strokes of a given character have two pairs of intersection points with the line segment.

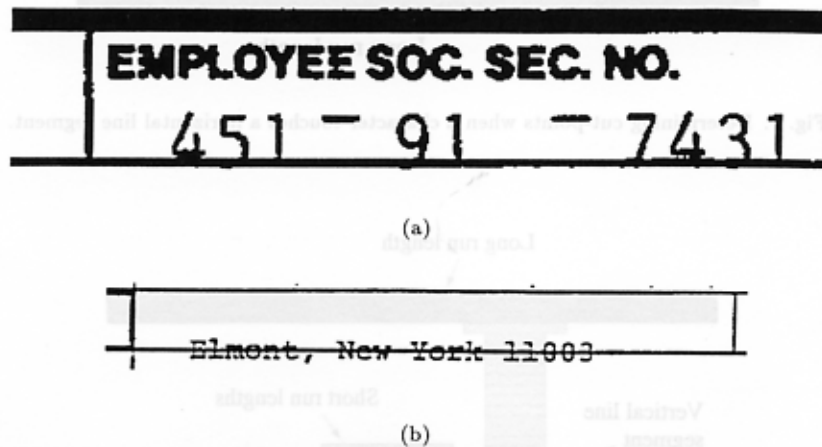


Fig. 6. Two cases of text-line contact: (a) characters touch a line segment, and (b) characters overlap a line segment.

In Sec. 4 we described the extraction of horizontal and vertical line segments from large connected components. Based on the information of the horizontal and vertical line segments, it is easy to design a method for separating touching text.

If a short run-length region does not include a vertical line segment, then it is considered as text that touches a horizontal line segment. An example where a character *E* is touching a horizontal line segment is given in Fig. 7. Light gray lines mean short run-length, and dark gray lines mean long run-length. Two ends of the short run-length that touches a long run-length are determined as the cut-point used to separate this character. For the overlapping case, we can use the same approach to decide the cut-points.

If a character touches a vertical line segment, the line segment and the character form a short run-length region (Fig. 8). As mentioned in Sec. 4, we can extract the vertical line segment from this short run-length region. So the cut-points A, B, C, and D in Fig. 11 can be decided as the intersection points between the boundary of the vertical line segment and four short run-lengths.

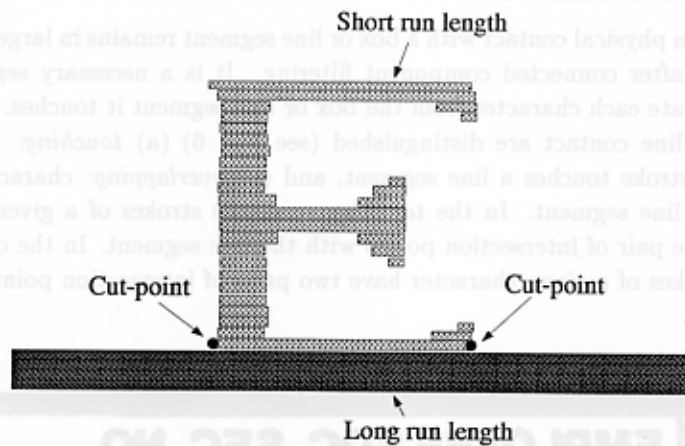


Fig. 7. Determining cut-points when a character touches a horizontal line segment.

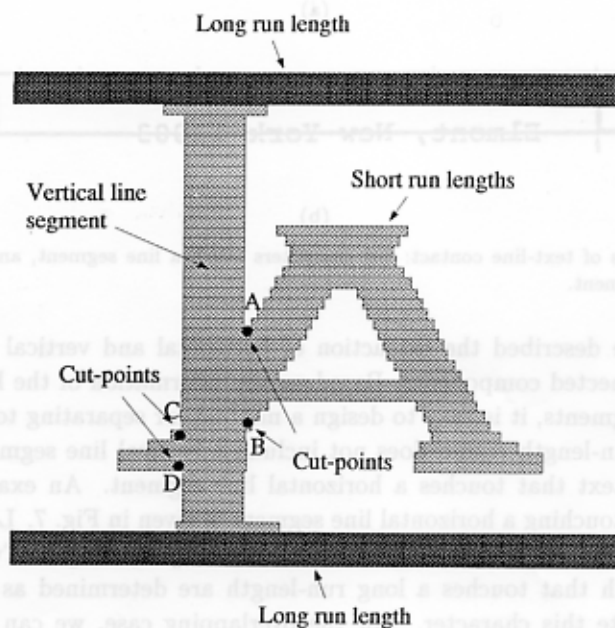


Fig. 8. Determining cut-points when a character overlaps a vertical line segment.

6. CHARACTER SPLICING

When a character is separated from a line segment that it touches, a part of the stroke will be missed after the line segment is removed. It may be necessary to

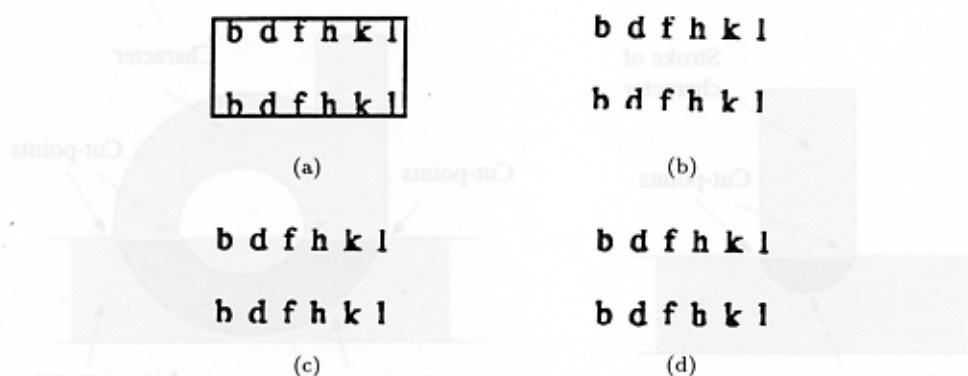


Fig. 9. Line removal and character patching: (a) an example of lines touching characters, (b) characters remaining after line removal, (c) result of arc patching, and (d) result of binding arc patching.

rejoin the missing parts of a character for the purpose of character recognition. The patching method should be different for the two cases where (i) characters touch lines and (ii) characters overlap lines, as discussed below.

Characters touch lines

An example of characters that touch lines is shown in Fig. 9(a). There are two subcases: (i) the character has only one pair of cut-points intersecting the line segment, and (ii) the character has more than one pair of cut-points intersecting the line segment.

(i) For some simple characters such as "i" and "l", it is easy to patch, e.g., by taking off part of the line segment which is located on the direction of elongation of the character stroke. But for some other characters such as "c" and "e", it is difficult to patch because of stroke curvature. For this we use *arc patching*: the character is patched by forming an arc between two cut-points, where the radius of the arc is chosen to be equal to one half of the average character height; see Fig. 11(a).

(ii) When there are more than one pair of cut-points intersecting a line segment, such as "b", "d", "h", and "k" in the bottom line of Fig. 9(a), the arc patching method mentioned above can be used to patch characters "h" and "k" as shown in Fig. 9(c). However, it is insufficient for patching characters "b" and "d". For such characters, another patching method, called *binding arc patching*, is used; see Fig. 11(b). Two arcs connect two pairs of cut-points without crossing the arcs. The result of binding arc patching on Fig. 9(a) is shown in Fig. 9(d); it can be seen that there are mistakes on characters "h" and "k".

One problem is the choice of patching method for a particular character in a particular form. The best way to patch characters is to set a feedback path between the form understanding system and the form analysis system. If a touching character has a choice of more than method of patching, then an arbitrary choice

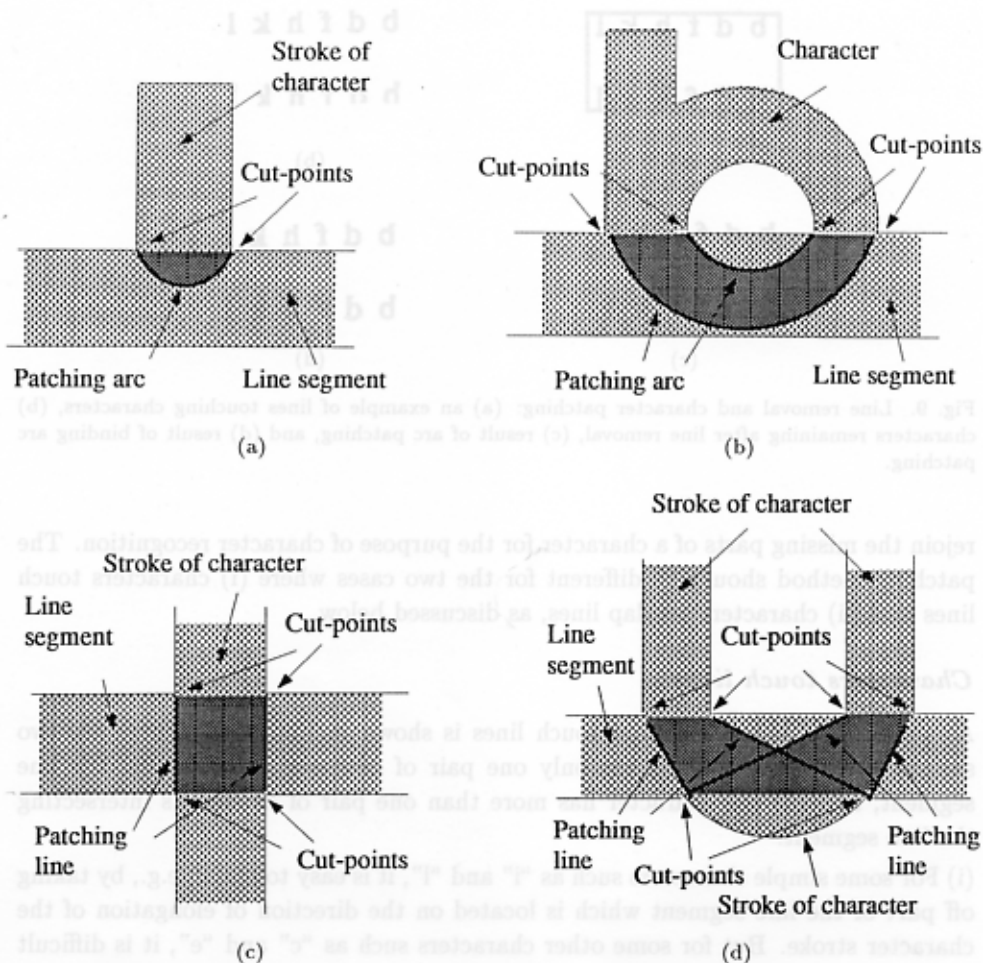


Fig. 10. Patching methods: (a) arc patching, (b) binding arc patching, (c) quadrilateral patching, and (d) quadrilateral patching.

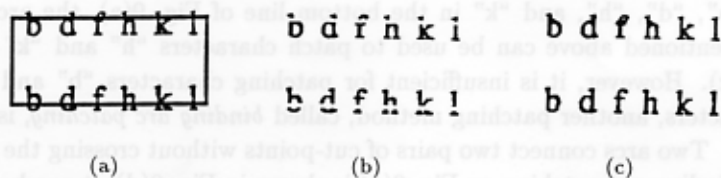


Fig. 11. Overlap repair: the lines in (a) are removed to yield broken characters in (b) which are patched to yield whole characters in (c).

is made and sent to the recognition system. The recognition system will send back a requirement to repatch this character if it is unable to recognize it. The form analysis system will patch this character again by another method.

Such an experiment was performed on Fig. 9(a). At first we chose the arc patching method to patch broken characters, and obtained the result shown in Fig. 9(c). Then the patched characters were sent to a character recognizer; the recognizer developed in our laboratory¹⁶ extracts structural features from characters and classifies them by a Bayes classifier. As a result all characters could be recognized exactly except the "b" and the "d" in the bottom line. The binding arc patching method was used to repatch these characters again, and the patched characters were correctly recognized.

Characters overlap lines

An example of lines overlapping characters is shown in Fig. 10(a). These are again two subcases. (i) The stroke of a character has two pairs of cut-points on both borders of a line segment respectively. The character can be patched by connecting corresponding cut-points and forming a quadrilateral for patching as shown in Fig. 11(c). (ii) A stroke of character has two pairs of cutpoints on one border of a line segment and another a pair of cut-points on another border of this line segment as shown in Fig. 11(d). For patching characters, quadrilaterals are formed by connecting one pair of cut-points on one side of line segment to both two pairs of cut-points on another side of this line segment as shown in Fig. 10(d). Characters separated from Fig. 10(a) are shown in Fig. 10(b), and Fig. 10(c) gives the patching result.

RELATIONSHIP	BIRTHDATE
Wife	5/7/56

(a)

RELATIONSHIP	BIRTHDATE
Wife	5/7/56

(b)

RELATIONSHIP	BIRTHDATE
Wife	5/7/56

(c)

Fig. 12. Examples of patching: (a) enlarged part of form image of Fig. 14(a), (b) broken characters after removing touching line segments, and (c) patched characters.

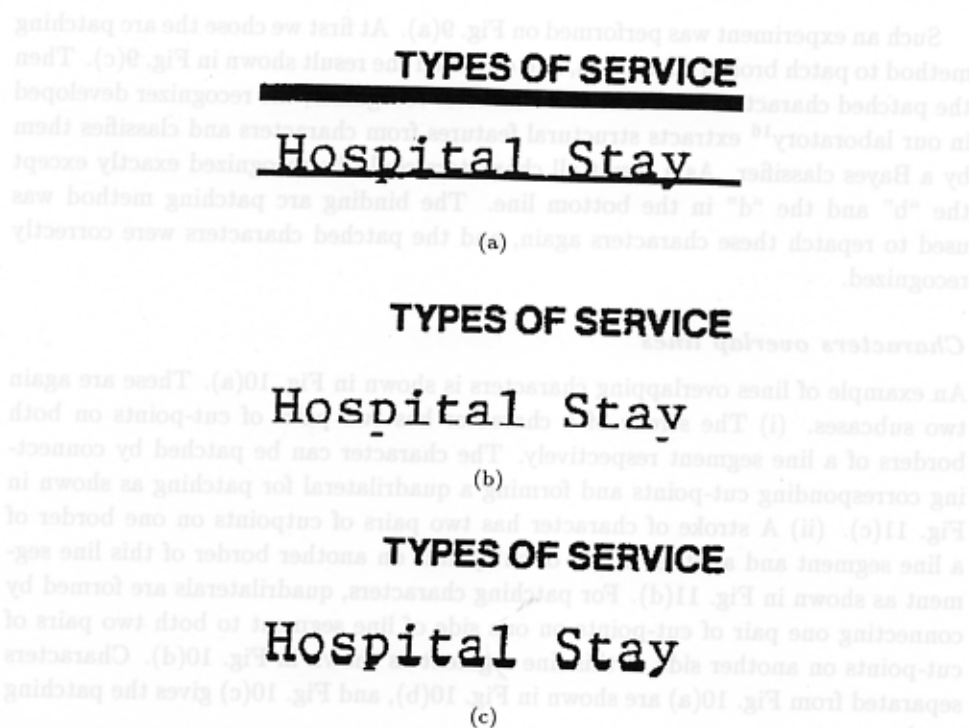


Fig. 13. Examples of patching: (a) enlarged part of form image of Fig. 14(a), (b) broken characters after removing touching line segments, and (c) patched characters.

Results of patching the form image of Fig. 14(a) are given in Figs. 12 and 13. In each, (a) shows an enlarged part of an original image which includes characters touching line segments, (b) shows broken characters after removing touching line segments, and (c) gives patching results.

7. EXPERIMENTAL RESULTS

Results with several forms, primarily pertaining to insurance and health-care, and containing both typed and handwritten entries are discussed here. Several test forms were digitized using a MICROTEK MSF-300GS Image Scanner at two resolutions of 200 and 300 ppi respectively. The system has been implemented in the C programming language under UNIX on a SUN Sparcstation 1(12.5 MIPS) computer.

Forms with typed entries

The form shown in Fig. 14(a) contains typed and handwritten entries. This image was digitized at 300 ppi resolution. The text map obtained by processing this image with the form analysis system is shown in Fig. 14(b). All boxes and line segments were removed correctly. The touching characters were separated exactly and spliced. The user time was 74 seconds, and the system time was 18 seconds.

[illegible]

(a)

Fig. 14. Example form tested in this research: (a) insurance form containing preprinted text, boxes, lines, typed and handwritten entries; the boxes have closed boundaries; (b) complete text map segmented from (a).

The text map corresponding to the form of Fig. 15(a) is given in (b). This form image was digitized at 200 ppi resolution. The user time was 23 seconds, and system time was 5 seconds. No significant errors resulted, with the exception of small check blocks being recognized incorrectly as characters. This is due to the inability of the system to distinguish between characters and small symbols.

This form analysis system was also run on thirteen form images which included a bank check example, an airline ticket example, a customer order form example, and a tax return form example. The results are quite satisfactory. All boxes and

EASTMAN KODAK COMPANY ELECTIVE BENEFIT ACCOUNT (EBA) • MEDICAL EXPENSE PLAN (MEP) • DEPENDENT CARE EXPENSE PLAN (DCEP)	 Kodak EBA P.O. Box 29 Rochdale, New York 14601 (716) 283-6510—no Rochester (800) 462-6725—for other calls within New York State (800) 628-6212—for calls from the remainder of the country
---	--

A-433005-FP

CLAIM (Please Print)

Part 1 EMPLOYEE INFORMATION

EMPLOYEE NAME (Last and First)	Vella, Henry	KODAK ID NO.	681-6345	EMPLOYEE AND SPOUSE NO.	651-91-7633
EMPLOYEE ADDRESS Number	341	Street	Agyle St.	City	Rochester NY Zip Code
EMPLOYEE TELEPHONE NO.	Home - 1 - (716) 691-6752		Work - 1 - (716) 651-8760		

Part 2 MEP EXPENSES

NITEM'S FULL NAME	RELATIONSHIP	BIRTHDATE	DATES OF SERVICE		TYPES OF SERVICE	CLAIM AMOUNT
			FROM	TO		
Belen Wells	wife	5/7/56	4/25	4/28	Hospital Stay	167.50

Part 3 DCEP EXPENSES

DEPENDENT'S FULL NAME	DATES OF SERVICE		TYPES OF SERVICE	CLAIM AMOUNT
	FROM	TO		
TOTAL MEP AND DCEP EXPENSES				167.50
				TOTAL MUST BE AT LEAST \$75.00

EMPLOYEE'S CERTIFICATION

I certify that expenses submitted under this claim:

- were for medical care as defined in MEP for me, my spouse, or my dependents, and/or for dependent care as defined in DCEP provided by an eligible provider;
- have not been and will not be otherwise reimbursed;
- in the case of DCEP expense are to be jointly employed;
- are for care that DCEP expenses accounted under this claim;
- ensure my spouse to be jointly employed or to seek work or are for services provided during months my spouse was an eligible disabled spouse or a full-time student;
- when combined with dependent care expenses reimbursed previously this year under DCEP do not exceed the lower of my or my spouse's earned income.

ANY PERSON WHO KNOWINGLY AND WITH INTENT TO INJURE, DEFRAUD, OR DECEIVE ANY INSURANCE COMPANY, FILES A STATEMENT OF CLAIM CONTAINING ANY FALSE, INCOMPLETE OR MISLEADING INFORMATION MAY BE GUILTY OF A CRIMINAL ACT PUNISHABLE UNDER LAW.

EMPLOYEE SIGNATURE _____ *Henry Wells*

DATE **5/12/88**

DO YOU...
☐ Attach EOBs or bills for MEP?
☐ Attach a receipt from the provider for DCEP?
☐ Keep copies for your records?
 Sign and date the form?

Employee's Copy

(a)

Fig. 14. Example form tested in this research: (a) insurance form containing preprinted text, boxes, lines, typed and handwritten entries; the boxes have closed boundaries; (b) complete text map segmented from (a).

The text map corresponding to the form of Fig. 15(a) is given in (b). This form image was digitized at 200 ppi resolution. The user time was 23 seconds, and system time was 5 seconds. No significant errors resulted, with the exception of small check blocks being recognized incorrectly as characters. This is due to the inability of the system to distinguish between characters and small symbols.

This form analysis system was also run on thirteen form images which included a bank check example, an airline ticket example, a customer order form example, and a tax return form example. The results are quite satisfactory. All boxes and

[illegible]

(a)

Fig. 15. Example form tested: (a) form containing preprinted text, boxes, lines, typed entries; the boxes do not have closed boundaries; (b) complete text map segmented from (a).

The input binary form image is first segmented into a partition of individual characters and a partition of boxes and line segments is obtained by connected component generation and filtering. Horizontal and vertical line segments are extracted from the partition of boxes and line segments. Key-points are extracted from the intersections of horizontal and vertical line segments as well as the ends of horizontal and vertical line segments. Based on key-points and their linkages, the structure of boxes are obtained and expressed as sequences of key-points. Since the horizontal and vertical line segments are extracted, text that touches the horizontal and vertical line segments can be separated from the line segments. After separation, the broken characters are spliced. Finally the system yields a complete text map and a list of all boxes and line segments.

Problems needing further investigation include cases where large characters are broken, and line segments are missed. Forms that have curved lines or have a reverse background also need to be studied.

5. K. Y. Wong, R. G. Casey, and F. M. Wahl, "Document analysis system", *IBM J. Res. Develop.* **26**, 6 (1982) 647-656.
6. I. Masuda, N. Hagita, T. Akiyama, T. Takahashi, and S. Naito, "Approach to smart document reader system", in *Proc. IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*, 1985, pp. 550-557.
7. L. A. Fletcher and R. Kasturi, "A robust algorithm for text string separation from mixed text/graphics images", *IEEE Trans. PAMI* **10**, 6 (1988) 910-918.
8. Y. Tsuji, "Document image analysis for generating syntactic structure description", in *Proc. Ninth Int. Conf. on Pattern Recognition*, 1988, pp. 744-747.
9. G. Nagy, S. C. Seth, and S. D. Stoddard, "Document analysis with an expert system", in *Proc. Pattern Recognition in Practice II*, Amsterdam, June 1985.
10. K. Inagaki, T. Kato, T. Hiroshima, and T. Sakai, "MACSYM: A hierarchical parallel image processing system for event-driven pattern understanding of documents", *Pattern Recogn.* **17**, 1 (1984) 85-108.
11. F. Esposito, D. Malerba, G. Semeraro, E. Annese, and G. Scafuro, "An experimental page layout recognition system for office document automatic classification: An integrated approach for inductive generalization", in *Proc. Tenth Int. Conf. on Pattern Recognition*, 1990, pp. 557-562.
12. J. L. Fisher, S. C. Hinds, and D. P. D'Amato, "A rule-based system for document image segmentation", in *Proc. Tenth Int. Conf. on Pattern Recognition*, 1990, pp. 567-572.
13. R. G. Casey and D. R. Ferguson, "Intelligent form processing", *IBM Syst. J.* **29**, 3 (1990) 435-450.
14. S. N. Srihari and G. W. Zack, "Document image analysis", in *Proc. Eighth Int. Conf. on Pattern Recognition*, 1986, pp. 434-436.
15. T. Pavlidis, *Algorithms for Graphics and Image Processing*, Computer Science Press, 1982.
16. D.-S. Lee, S. W. Lam, and S. N. Srihari, "A structural approach to recognize hand-printed and degraded machine-printed characters", in *Proc. Int. Assoc. for Pattern Recognition Workshop on Syntactic and Structural Pattern Recognition*, 1990, pp. 256-272.

Received June 1992.



Sargur N. Srihari, Professor of Computer Science, is the director of CEDAR, Center of Excellence for Document Analysis and Recognition, located at the State University of New York at Buffalo.

Dr. Srihari is the author of over 125 technical papers in the science and engineering of reading machines, artificial intelligence and pattern recognition. He is co-author of several US patents and a book on *Computer Text Recognition and Error Correction*.

Dr. Srihari has been an Assistant, Associate and full Professor at the State University of New York at Buffalo since 1978.

Dr. Srihari received his B.Sc. degree in physics and mathematics from Bangalore University in 1967, B.Eng. degree in electrical communication engineering from the Indian Institute of Science, Bangalore in 1970, and the M.S. and Ph.D. in computer and information science from the Ohio State University, Columbus in 1971 and 1976, respectively. He is a Fellow of the Institution of Electrical and Electronics Engineers, a member of the American Association for Artificial Intelligence, Pattern Recognition Society and the Association for Computing Machinery. He received a New York State/United University Professions Excellence Award in 1991.



Dacheng Wang is a research scientist in the Center of Excellence for Document Analysis and Recognition, State University of New York at Buffalo, USA. His main research interests include document image analysis and understanding, handwritten character recognition.

Wang received his B.S. degree from the Chinese People's Liberation Army Military Telecommunication Engineering Institute in China. He had been a research fellow in Kyoto University, Japan for two years.