

# Line Separation for Complex Document Images Using Fuzzy Runlength

Zhixin Shi and Venu Govindaraju

Center of Excellence for Document Analysis and Recognition(CEDAR)  
State University of New York at Buffalo, Buffalo, NY 14228, U.S.A.

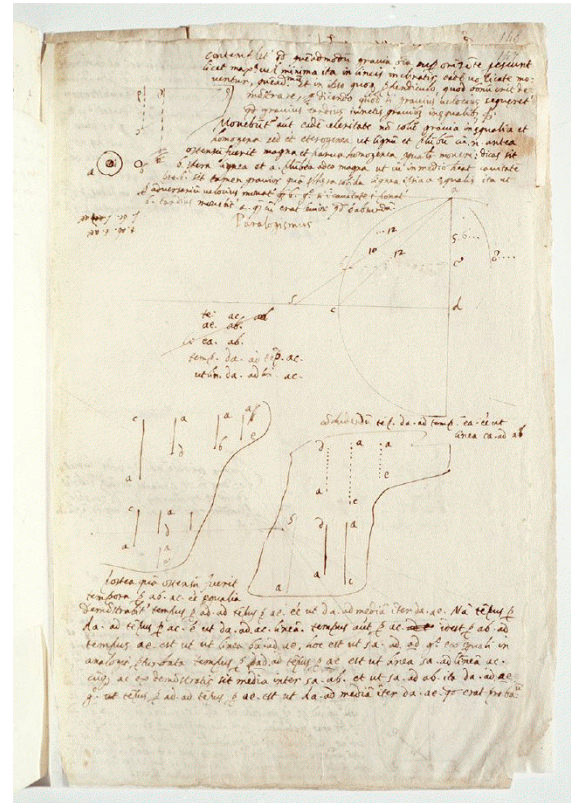
## Abstract

*A new text line location and separation algorithm for complex handwritten documents is proposed. The algorithm is based on the application of a fuzzy directional runlength. The proposed technique was tested on a variety of complex handwritten document images including postal parcel images and historical handwritten documents such as Newton's and Galileo's manuscripts. A preliminary testing showed a successful rate of 93% of the test set.*

## 1. Introduction

The text line location and separation of a document page is an important task for document analysis and optical character recognition (OCR) applications. Line separation algorithms first locate the text lines and then segment the text line in their original logical order. A wide variety of methods have been proposed in the literature. The Projection Profile technique[2, 6, 3] detects the text lines by creating a histogram at each possible location. The methods using Hough Transform are theoretically identical to the methods using projection profiles[1]. Hough Transform is usually used for locating skewed text lines. It is applied on a set of selected angles along each angle straight lines are determined with a measurement for the fit. The best fit for the lines gives the skew angle and the location of the lines. Hough Transform can be applied on all black pixels[7], on reduced data from horizontal and vertical run-length computations[5] or only on the bottom pixels of connected components[4]. Another method uses nearest neighbor clustering of connected components found in [10].

However, most of the proposed approaches are designed mainly for machine printed documents. They are usually able to deal with small skew angles within  $\pm 15^\circ$ , failing to manage cases of documents that may exceed this limit. Moreover, some of them entail high computational cost, especially in the case where the Hough transform is used. Also, certain approaches are font, column, graphics or border dependent. There are very few methods proposed



**Figure 1. Historical document image, Galileo's manuscript. Complexity is in the low image quality due to aging of the paper document, handwriting with touching between lines and mixed with graphics, etc.**

to handle handwritten or mixed documents[9, 8]. Some of them are designed for specific applications. Most of these methods can not handle handwritten documents very well. This is mainly because of the complexity of handwritten documents including variations in writing styles, sizes of characters and touching or connection of characters, words and text lines.

In this paper, we propose a novel text line location and

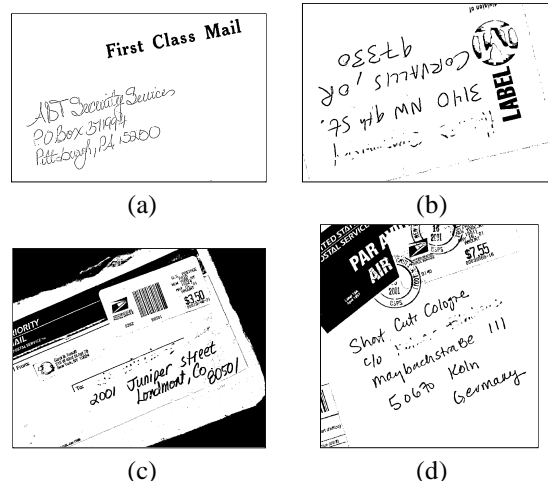
separation algorithm for complex document. The method is able to detect the per-dominant directions and locations of the text lines and separates the lines in the document image. The text can be handwritten, machine printed or mixed. The main advantage of the proposed method is the ability to extract the text line information from document mixed with non-text elements such as graphics, bar-codes and forms, see examples in Figures 1 and 2. The method first transforms the content of a document page into components representing text words, line or graphics areas. Then the extracted components are ranked as being possible types of text or non-text. Finally the per-dominant direction of the text in the image is estimated using minimum squared distance optimization and component grouping method, which is followed by separation of the text lines.

In Section 2 we will discuss the complexity of handwritten and mixed documents. The basic principle of our approach will be described. In section 3 the steps in locating text lines will be described in detail. Special care is taken in splitting touching lines. In Section 4 we will present experiment results and Section 5 gives a final conclusion.

## 2. Background

As indicated in [10], without recognizing character and/or understanding the document, human are usually be able to determine the locations and orientation of text lines in document images by collecting some regularly aligned symbols, mostly text words or characters. A symbol line is defined as a group of regularly aligned symbols that are adjacent, relatively close to each other and through which a straight line can be drawn. Connected components are chosen as such symbols in [10]. Foreground pixels can also be chosen as such symbols. For example in methods using projection profile or Hough transform, pixels are used to determine the text line directions along which the most favorable profiles or straight lines can be drawn. But in some complicated document images such as the images in Figure 1 and 2, obviously not all of the foreground pixels can contribute in making the symbols to determine the document orientations. The previous methods fail on some of these images. For example images (b), (c) and (d) in Figure 2 can easily mislead the projection profile or Hough transform because of the multiple directions of text lines or the noise, the big black border and graphics. In image (a) in Figure 2, there are both too small and too big connected components from the broken text and connected text crossing lines. They make the methods using connected component approaches also hard to find the correct directions.

For the purpose of document recognition/understanding, we are mostly interested in determining the document texts



**Figure 2. Complex images :** (a) Connected components in the text lines are too big and crossing lines. (b) Multiple text lines. (c) and (d) Mixed images of big black foreground, printed text, handwritten text, noise and other graphics.

regardless of they may be buried in other document object such as graphics. Without recognition, to separate the texts from the other graphic elements is a difficult issue. Especially in handwritten documents, we very often see touching characters not only between adjacent neighboring characters, but also between characters in different text line. On the other hand, due to scanning quality or imperfection in binarization processing, we also see many characters in text are broken to smaller pieces. To correctly identify text lines without recognition, we have to combine grouping of the broken segments and separating of the touching or connected characters between lines together.

From our analysis of the complex document images we found some general properties of text lines and other graphics. We noticed that for both handwritten and printed text, the relative distances between characters in same lines are generally smaller than the distances between text lines. Human identification of the text lines utilizes these differences efficiently. We can also tolerate touching or connection between text lines by finding a background path between each pair of text lines. The touching or connections between text lines are usually made by oversized characters or characters with long descenders running through the neighboring lines. Therefore if we could ignore some bridging strokes crossing lines, we may see the background paths between lines clearly. As for graphics, they are either in relatively irregular shapes as isolated connected components or a group of connected components occupying a bigger area than any

usual text line. Examples are bar-codes, stamps, pictures and pepper noises.

Based on these observations, we decide to use the background runs to build a type of separators. They should separate different text lines by running through the connecting strokes between lines. They should also be able to group the neighboring characters in same text lines together.

Background runlength has being used in the literature for page segmentations and skew detections[6]. Background regions(white streams) are built from adjacent background runs and those wide white streams are used in estimation of skew angles. For page block segmentation, the assumption that the text lines and blocks are well separated by white background is required. For tolerating noise and run-away black strokes, runlength smearing[7] is applied as an image processing. The favorite runlength such as foreground runs are created by skipping small runs in background color. The expected result from this process is that the most of the foreground text characters are grouped together. The text lines and text blocks are extracted by using a connected component analysis approach. The method again works well for printed documents with mostly text. It will fail on documents with touching or connection between text lines and text blocks.

We may use the runlength smearing approach for building our background runs. “Smearing” is as same as skipping or ignoring some foreground pixels. Hopefully it will break the touching between text lines. But setting up the threshold for skipping the small foreground runs is difficult. If we set the value too small, we may not be able to break the crossing strokes connecting the text lines. If we set the value too big, then it may be bigger than the stroke width of most of the characters and end up erasing the text areas.

To solve the problem we designed a new kind of runlength – fuzzy runlength. we *trace* a background run starting from a background pixel along two directions, to its left and right(this is for horizontal runs, otherwise the up and down directions for vertical runs). On the way of the tracing we skip some foreground pixels. When the accumulated number of skipped pixels exceeds a pre-set threshold, we stop the tracing and do the same along the other direction. The total number of the traced positions is the length of the run associating to the position where we start the tracing. Intuitively, the fuzzy runlength at a pixel is how far we can see from standing at the pixel along horizontal(or vertical) direction. Like a human standing in a forest looking for a path out of the forest, the length that he can see along a direction may not be the distance from where he is standing to the first tree in front of him. Rather he may be able to “see through” a few trees to get a longer view. But he may not be able to see through too many trees.

At the end of the tracing process, what we get is a two dimensional matrix in the size of the original binary im-

age. Each entry of the matrix is the fuzzy runlength of the pixel in its position. We take this matrix as an image and binarize it by using another pre-set threshold. A pixel is set to be background if its pixel value is bigger than the threshold, otherwise we set it as foreground. Most of the background pixels are originally background pixels in the original document image. The foreground consists of connected components made of blurred blocks of original foreground pixels. As we expected, fuzzy runs break the touching text lines and group the characters in the same lines to their close neighbors. The text lines, mostly text words are among the foreground connected components.

To determine the orientation of the text lines, we only need to identify some of the foreground components covering the text lines. For this purpose we will apply a size constraint on the components. This size constraint could be either pre-set value for average height of a usual text line or an estimated height from connected components in the original document image.

In next section we will present the algorithms of constructing the fuzzy runlength and estimation of text line locations and orientations using a selected set of above foreground components.

### 3. Text lines detection

#### 3.1. Build fuzzy runlength

We assume that the input document images are binary images with foreground color black and background color white. We generate the horizontal fuzzy runlength by scanning each row of the image twice, from left to right and then from right to left.

1.  $In\_img \leftarrow$  the input binary image.
2.  $Fuzzy\_img \leftarrow$  out put buffer of same size as the input image, initialized to 0 for holding the fuzzy runlengths.
3. Scan each row of the input image from left to right:
  - (a) Initialize a FIFO queue **blkQ** for holding black pixels which can be seen through from the current pixel to its left.
  - (b) For the current position  $j$ , if the current pixel is black, put the current position  $j$  on top of the queue **blkQ**.
    - \* If the number of black pixel in **blkQ** is bigger than **MaxBlockCnt**, assign the fuzzy run image at the current position the value of  $j - \mathbf{blkQ}[\mathit{bottom}]$ . Then pop out the bottom element from **blkQ**;
    - \* Else assign the fuzzy run image at the current position the value of the fuzzy run image at previous position + 1.

4. Similarly we scan the row from right to left.

In the above algorithm **MaxBlockCnt** is a threshold value we have to decide before scanning the image for calculating the fuzzy runs. It can be either set as a fixed value for an application running similar set of images or determined at run time dynamically. Since the value of **MaxBlockCnt** is for determining how many black pixels we want to *see-through* in building the fuzzy background runs, it can actually be set as a large value, as far as it is not too large so that we can see through all the black pixels in a sizable word. For example, in each word if we assume there are at most two descendent characters such as g or y, and the average stroke width is 5 pixels. We want the see-through pass 4 strokes then we can set **MaxBlockCnt** to be a little more than 20, say, 25.

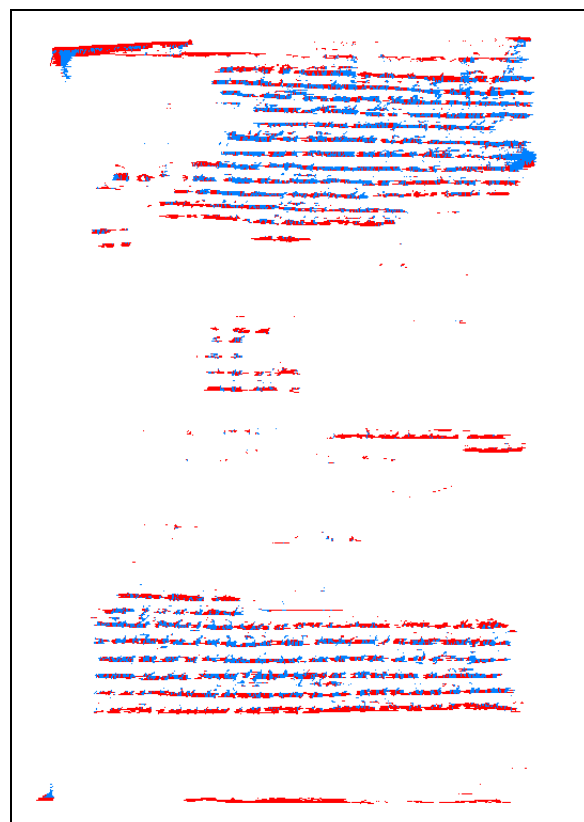
### 3.2. Locations of possible text lines

After above scan we get a buffer *Fuzzy\_img* for holding fuzzy runlengths for each pixel. We take the buffer as an image and binarize it to two values. The threshold value for the binarization can be a pre-set value or can be determined dynamically. The fuzzy runlength at each pixel represents the maximal extend of the background along the horizontal direction through the pixel position. Considering the runs inside a text word being shorter than the runs in between the text lines, we can choose the threshold value as bigger than the estimated maximal runlength inside the words. In our experiment we empirically determined the threshold based on the estimations of average character size, average distance between characters and average stroke width. These values can all be calculated dynamically too.

The binarized fuzzy runlength image *Fuzzy\_img* in Figure 4 consists of connected components each as a pattern represents either part of a text line(made of one or a few words) or other graphic elements. For a document with text line skew angle between  $\pm 45^\circ$  the fuzzy runlengths can push the text line patterns clearly stand out.

To identify a connected component as a pattern of a text line, we have to give a clear definition of text line pattern. We define a component as a pattern of a text line if it satisfies the following conditions:

1. The height of the component should be near the estimated text height. This height can not be the height of the bounding box of the component. To calculate the height we evenly divide the component to many small pieces horizontally and calculate the vertical extend from each one of these pieces. The average values of these vertical extends is the estimated height of the component, see Figure 6.
2. the length of the component should be long enough, for example it should be at least twice as long as the cal-



**Figure 3. Fuzzy runlength shows very good patterns for text lines. The touching lines are well separated.**

culated height of it. For the length we simply use the length of the bounding box.

3. We may also add a condition on its density to filter out noise. This is optional. The identified text line pattern are shown in Figure 5.

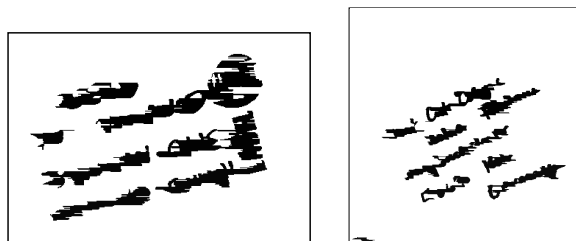
### 3.3. Detection of skew angle of text areas

Each text line pattern as a connected component is a long rather than wide strap of black pixels. It is like a blurred image of a text line or word. Since the descenders and ascenders are striped off, its orientation along the longitude represents the orientation of the underlying text line or word. Therefore we first compute the orientation for each of the components.

For each text line pattern we first evenly divide it horizontally into  $n$  pieces, see Figure 6. The distance  $d$  between the adjacent dividing lines is fixed. When  $d = 1$  we simply use all the columns in its bounding box as dividing lines. Choosing a bigger value for  $d$  is a matter of efficiency. On



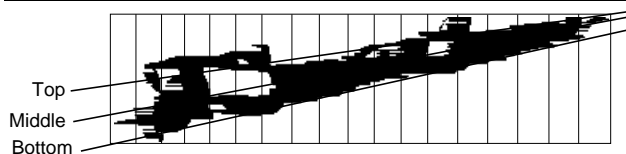
**Figure 4. Binarized fuzzy runlength images. Connected components show patterns for text lines. Line patterns for text lines stand out from other graphics clearly. They can be identified using a general shape constraint.**



**Figure 5. Identified connected components as text line patterns.**

each dividing line we take three points, the top most black pixel, the lowest black pixel and the middle point between these two pixels. We name them *top*, *middle* and *bottom* points. We then use minimum squared distance method to get three best fit straight lines for all the top points, all the middle point and all the bottom points. Compare the fits (the minimal squared distances) of the three lines we choose the best among them. In our example in Figure 6 the bottom line is our best choice. The orientation of the line pattern is defined as the orientation of the chosen straight line.

There are several other ideas for computing the orientation. One of them is to get the estimated convex hull of a



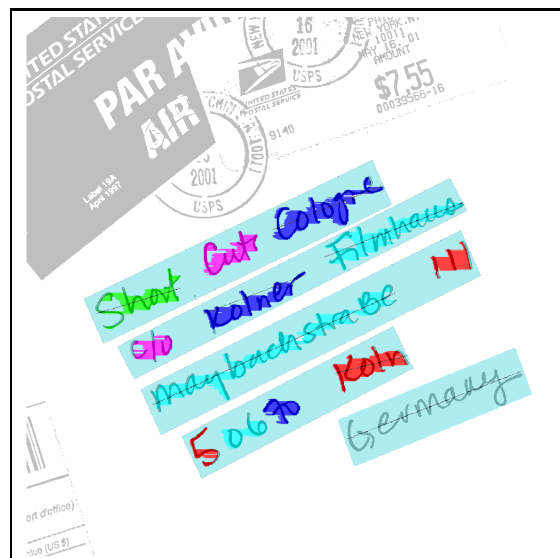
**Figure 6. Compute a best fit orientation for a component**

component first, then use the convex hull to compute the orientation.

### 3.4. Grouping text line patterns

Now that we have detected text line patterns each is a cover of a text line or part of a text line. We then use the orientation information for each text line patterns to group them into complete text lines. To do so, we follow a heuristic approach.

Since the orientations of relatively small pieces are not accurate enough, we first rank all the available information base on not only the orientation of each component, but also other information including the sizes and the shapes of the components. For example the longer components should give more reliable orientation estimations. Then we start with the most reliable component. Each component will be grouped into a group which has closet orientation and location match so to form a text line. A new group will be started if a component can not be put into an existing group due to the constraints such as height of a line formed by the components in a group should not be increased significantly. This grouping procedure continue until all the detected line components are exhausted, see Figure 7.



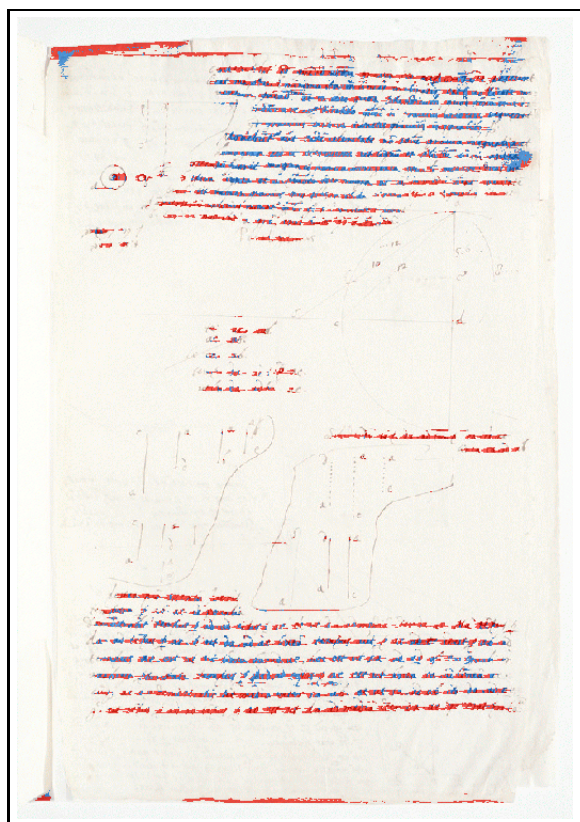
**Figure 7. Grouped line components give the locations of the text lines.**

## 4. Experiment

To test our method, we used a set of 1864 USPS parcel images. These images are cropped and binarized from

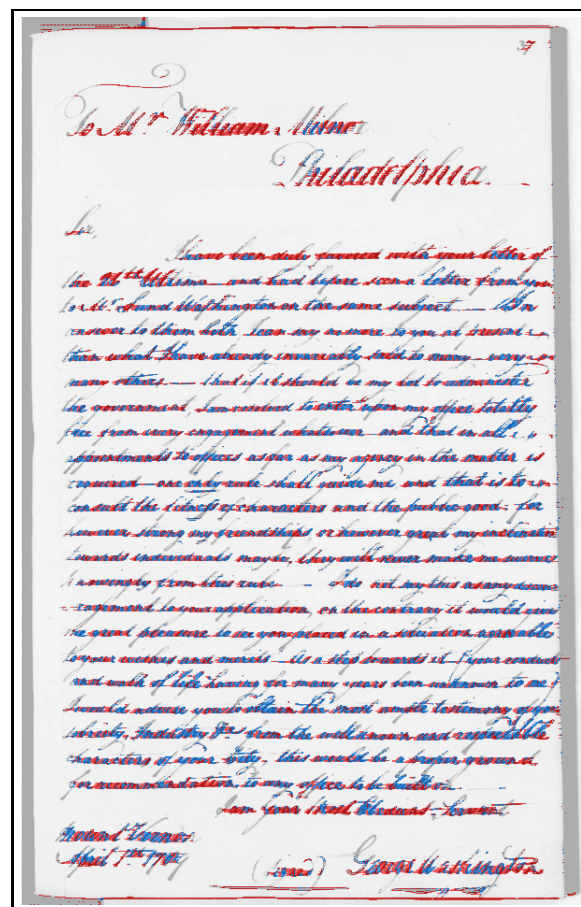


the original parcel address images automatically. The original scan resolution was 300dpi. For locating text lines, we don't need such high resolution. We down sampled them to 1/4 of the size before running our programs. The correct rate of our algorithm is 93%. The similar is also done on a small set of historical handwritten documents manually. These images include Newton's manuscripts, Galileo's manuscripts and Washington's manuscripts. Examples of the experiment shown in Figure 8 and 9.



**Figure 8. Extracted text line pattern superimposed on top of Galileo's manuscript. Grouped line components will give the locations of the text lines.**

Among the failure cases are the images failed from being using fixed parameters in our simple implementation. For example the threshold value for "running through" black pixels is too small for images with extremely large size characters with thing strokes. The dynamically determined such threshold value based on evaluations of character size and stroke width would help for the problem.



**Figure 9. Extracted text line pattern superimposed on top of Washington's manuscript.**

## 5. Conclusion

In this paper we present a novel method for complex document text line separation. The method uses a new concept of fuzzy runlength which imitates an extended running path through a pixel of a document. The fuzzy runlength can be use as an image processing tool to partition a complex document to texts in terms of text words or text lines, and to other graphic areas. Classification of these areas can lead us to getting information for document segmentations especially text line separations. Our simple experiment also showed the success of the method. Further research along this direction will be using the method on document segmentation and other content locations.

## References

- [1] A. Amin, S. Fischer, T. Parkinson, and R. Shiu. Fast algorithm for skew detection. *IS&T/SPIE Symposium on Electronic Imaging, San Jose, USA*, 1996.
- [2] H. Baird. The skew angle of printed documents. *Proc. SPSE 40th Conf. Symp. Hybrid Imaging Syst., Rochester, NY*, pages 21–24, 1991.
- [3] G. Ciardiello, G. Scafuro, M. Degrandi, M. Spada, and M. P. Roccotelli. An experimental system for office document handling and text recognition. *Proc 9th Int. Conf. on Pattern Recognition*, pages 739–743, 1988.
- [4] G. T. D.S. Le and H. Wechsler. Automated page orientation and skew angle detection for binary document image. *Pattern Recognition*, 27:1325–1344, 1994.
- [5] S. Hinds, J. Fisher, and D. D’Amato. Document skew detection method using run-length encoding and the hough transform. *Proceeding of 10th International Conference on Pattern Recognition*, pages 464–468, 1990.
- [6] T. Pavlidis and J. Zhou. Page segmentation by white streams. *Proc. 1st Int. Conf. Document Analysis and Recognition (ICDAR), Int. Assoc. Pattern Recognition*, pages 945–953, 1991.
- [7] S. Srihari and V. Govindaraju. Analysis of textual images using the hough transform. *Machine Vision and Applications*, 2:141–153, 1989.
- [8] A. H. W. Chin and A. Jennings. Skew detection in handwritten scripts. *Proc. IEEE on Speech and Image Technologies for Computing and Telecommunications*, pages 319–322, 1997.
- [9] B. Yu and A. Jain. A generic system for form dropout. *IEEE trans. On Pattern Analysis and Machine Intelligence*, 18(11):1127–1134, 1996.
- [10] B. Yu and A. Jain. A robust and fast skew detection algorithm for generic documents. *Pattern Recognition*, 9:1599–1629, 1996.