

DOCUMENT IMAGE ANALYSIS

Sargur N. Srihari

Department of Computer Science
State University of New York at Buffalo
Buffalo, New York 14260

Gregory W. Zack

Imaging Systems Laboratory
Xerox Webster Research Center
Webster, New York 14580

ABSTRACT

A document image is an optically scanned and digitized representation of a printed page. Typically it consists of fields of text, line drawings, half-tone pictures, and icons. Document image analysis is the problem of deriving and labeling such fields. This paper surveys several techniques that should prove useful for the task.

1. INTRODUCTION

A document image is a visual representation of a page from a journal, magazine, newspaper, etc. Typically, it consists of blocks of text, i.e., letters, words, and sentences that are interspersed with half tone pictures, line drawings, and symbolic icons. A digital document image is a two-dimensional array obtained by either optically scanning and raster digitizing hard copy or an electronic version created for a laser printer.

Document image analysis is the task of deriving a high level representation of the contents of a document image. As in the case of any image analysis task, the representation is application dependent. The structural elements of a document, i.e., columns, paragraphs, titles, figures, lines of text, symbols, etc., are generally laid out in rectangular blocks aligned parallel to the horizontal and vertical axes of the page. An editable description will include the following features:

- major components, e.g., text, line drawings, tables, graphs, photographs are identified,
- spatial relationships between components, e.g., reading order, are derivable,
- layout features, e.g., number of columns, margin justification, are identified,
- in case of drawings, objects are identified, e.g., straight lines, curves, solids and arrow heads,
- in case of text, roles are assigned, e.g., title, subtitle, caption; fonts are distinguished and characters/words are recognized.

Document analysis is the inverse of electronic document synthesis (Fig. 1). In the latter case, a high-level description is converted to a raster format whereas in the former a raster image is converted to a high-level description.

Document Analysis Components

The three phases of a document analysis process are shown in Fig. 2. Phase 1 consists of optical scanning, digitization, and binarization. Phase 2 consists of block segmentation and labeling. The result of phase 2 is a set of regions (possibly overlapping) labeled as text, graphics, half-tone, etc. Phase 3 consists of several parallel operations for processing text, graphics, and half-tone images. Certain

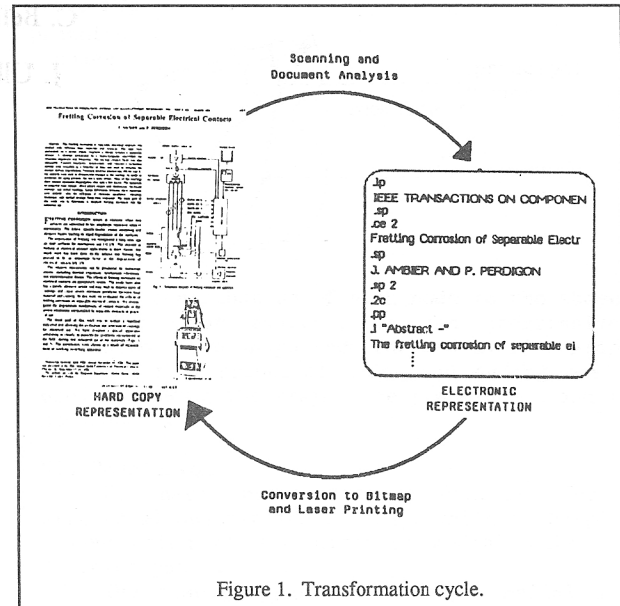


Figure 1. Transformation cycle.

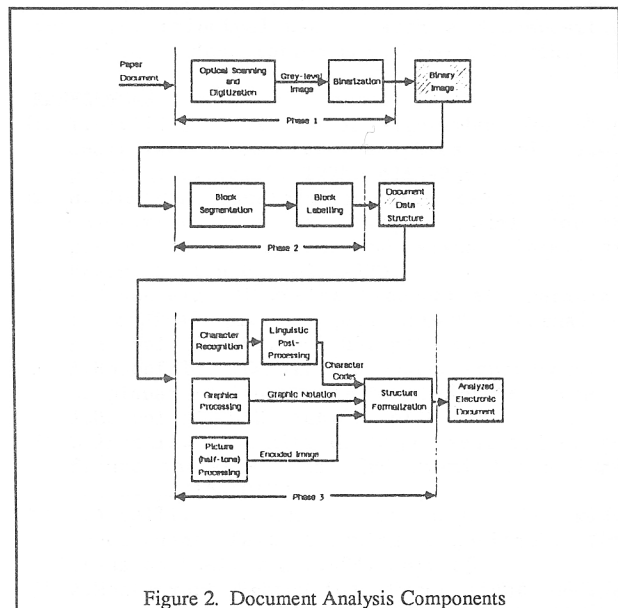


Figure 2. Document Analysis Components

document analysis problems may require feedback paths that are not shown here. The remainder of this paper discusses Phase 2 operations only.

2. BLOCK SEGMENTATION

Approaches for segmenting document image components may be divided into two camps. *Top-down*, or knowledge based, techniques proceed with an expectation of the nature of the document. They first divide the document into major regions which are further divided into sub-regions, etc. *Bottom-up*, or data-driven methods progressively refine the data by layered grouping operations. Although no system takes a pure approach, most can be identified as being aligned with one philosophy.

2.1 Top-Down Approaches

Approaches that may be described as top-down are smearing [1], projection profile cuts [2], and Hough transform detection [3].

Smearing

A run-length smearing algorithm (RLSA) operates on the image under which any two black pixels (1's) which are less than a certain threshold t apart are merged into a continuous stream of dark pixels. White pixels (0's) are left unchanged. The RLSA is first applied row-by-row and then column-by-column, yielding two distinct bit maps. The two results are then combined by applying a logical AND to each pixel location. The resulting RLSA image contains a smear wherever printed material appears on the original image.

Projection Profile Cuts

Based on the observation that printed pages are primarily made up of rectangular blocks, a page can be recursively cut into rectangular blocks. Thus the document is represented in the form of a tree of nested rectangular blocks. The application of cuts is based on the configuration of the pixels. A "local" peak detector is applied to horizontal and vertical "profiles" to detect local peaks (corresponding to thick black or white gaps) at which the cuts are placed; it is local in that the width is determined by the nesting level of the recursion, e.g., gaps between paragraphs are thicker than those between lines.

A complete hierarchical recursive tiling of a page can be represented as an x-y tree. It is obtained by using alternating horizontal or vertical partitions [4]. The root represents the entire page and can be considered either a single horizontal or vertical partition. The number of cuts placed at each level is variable, so the resulting tree is not necessarily binary. A node that represents a horizontal (vertical) partition has a set of descendants that represent a complete vertical (horizontal) tessellation of the parent node.

Hough Transform

A document typically contains several straight lines. Forms and tables contain solid lines. Line drawings, e.g., block diagrams, predominantly have straight lines. Columns of text are separated by straight rivers of white space. Text usually consists of several parallel textured thick lines. The Hough transform is a technique for detecting parametrically representable forms, e.g., straight lines, in noisy binary images. In the case of detecting straight lines the transformation used is $\rho = x \cos \theta + y \sin \theta$. The accumulator array resulting from applying the Hough transform to the "abstract" portion of the document image of Figure 1 is shown in Figure 3. In general the array can be checked for the particular orientation which has the maximum

number of transitions to and from a minimum value. Transitions corresponding to text are usually regular and uniform in width and thus are easy to identify.

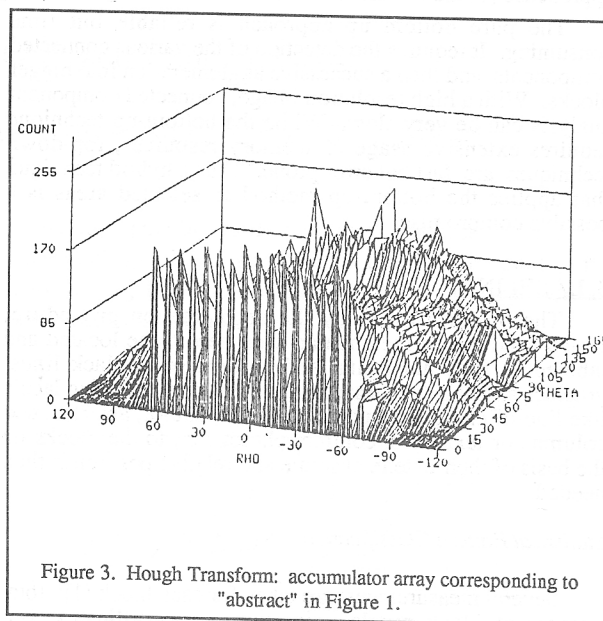


Figure 3. Hough Transform: accumulator array corresponding to "abstract" in Figure 1.

2.2 Bottom-Up Approaches

The image is first processed to determine the individual connected components. At the lowest level of analysis are individual characters and other large figures. In the case of text the characters are merged into words, words are merged into lines, lines into paragraphs and paragraphs into even larger blocks, if such a merging is possible [4]. It is usually necessary to determine whether a connected component is a part of: text, line drawing, or threshold region of a half-tone picture. Possible features for performing this classification are: size, branching structure, topology, and shape measures.

The geometrical complexity of an object can be measured by the number of linear components (branches) that the figure is composed of. If the complexity is high then the object is probably part of a line drawing, half-tone picture, or form. An algorithm can be designed to determine geometrical complexity of several components while the image is scanned raster fashion [6]. This is done by keeping track of the components in each line and their connectivity with respect to components in the previous line.

Topological properties of a component are useful in determining if it is a character. The Euler number of a binary image is defined as the number of components minus the number of holes. For a component which is a character the Euler number is one (1), zero (0), or negative one (-1).

2.3 Critical Evaluation of Approaches

Smearing is fast, but fails if text lines are skewed. The thresholds have to be set a priori; thus it is specialized to a certain class. Block size, eccentricity and the ratio of the number of black pixels in the original image to the number of pixels in the segmented block is not sufficient to distinguish between text/graphics, e.g., a dense field of random straight lines would be classified as text. The x-y tree approach also fails with skewed lines. In the case of a noisy document, the cuts may be erroneously placed. Once cuts are placed it is not

possible to recover, since cut-building is a generative process. The Hough transform and other histogram approaches [7] are invariant to skew.

The pure bottom up approach is reliable but time consuming. It requires the detection of the various connected components and then a successive agglomeration into bigger blocks. With a high resolution image, connected component analysis can be very slow. While the bottom-up technique requires extensive usage of memory resources, top-down techniques are fast but error prone. Thus a hybrid technique that applies the bottom-up method in selected areas is a possible compromise.

3. LABELING

The blocks formed by the segmentation procedure, which need not be rectangular, then need to be located and labeled as text, horizontal and vertical solid black lines, graphics and half-tone images. A method of associating labels like line-of-characters, paragraph, title, author, abstract, text columns or text lines, column, photo, etc., to the blocks on the basis of their extent, absolute and relative position is then needed.

Statistical Pattern Classification

Several measurements are taken for each block [1]: total number of black pixels, minimum x - y coordinates, x - y lengths, total number of black pixels, and number of horizontal white-black transitions in the original image block. Classifying each block is done by computing several features for each block from the above measurements and then using a linear pattern classifier.

A statistical method can be devised to determine whether a block contains only text. It involves computing several filter histograms of the test region. One example of a filter is a white-to-black transition separated from another white-to-black transition by n units along the horizontal direction. The histogram is obtained by counting the number of such transitions that occur for each value of n . The histograms for printed text tend to have several regular peaks while handwriting tends not to have such peaks. The histogram can be classified by treating it as a feature vector.

Rule-Based Approach

Labeling can be done by a rule-based expert system that uses a knowledge-base of layout and composition rules for specific classes of documents [4,8]. These rules can determine line-of-characters, paragraphs, columns, photos etc. using width, aspect ratio, number of cuts, number of siblings at a particular stage and other such related information.

4. CONTROL STRUCTURES

The coordination of several processes is of central importance in document image analysis. Methods based on production systems and blackboard communication [9,10] hold promise.

A production system that is organized with different levels of rules is given in [8]. The control flow is as follows: the image is segmented and data about various regions are obtained. This data includes intrinsic properties (shape, size, aspect ratio, etc.) as well as spatial relationships between regions. Knowledge rules pertain to intrinsic properties and

spatial relationships. Control rules decide what knowledge rules are to be executed and in what order, and thus act as focus-of-attention mechanisms to guide the search. Strategy rules supervise the entire search and classification process, and determine whether a consistent interpretation of the image has been obtained. If the data from the initial segmentation is insufficient, then the system invokes further processing to obtain more data.

There are several advantages to using production rules in this domain. Unlike natural scenes, documents are highly structured, and thus features can be very effectively formulated as production rules. It is easy to apply either an additional strategy to a region that is hard to interpret, or a retry process having modified parameters. Addition/modification of rules is a relatively simple process that does not disrupt the rest of the system.

REFERENCES

- [1] K. Y. Wong, R. G. Casey and F. M. Wahl, "Document analysis system", *IBM J. Res. Dev.*, 26 (6), 1982.
- [2] H. Zen and S. Ozawa, "Extraction of the fair document from mixed mode manuscript", *Proc. CVPR*, San Francisco, 1985, 544-549.
- [3] A. Rastogi and S. N. Srihari, "Recognizing textual blocks using the Hough transform", TR 86-01, Dept. of CS, SUNY/Buffalo, 1986.
- [4] G. Nagy, S. C. Seth, and S. D. Stoddard, "Document analysis with an expert system", *Proc. Pat. Recog. in Practice*, Amsterdam, 1985.
- [5] W. Doster, "Different states of a document's content on the way from the Gutenbergian world to the electronic world," *Proc. 7th ICPR*, 1984, 872-874.
- [6] T. Pavlidis, "A hybrid vectorization algorithm," *Proc. 7th ICPR*, 1984, 490-492.
- [7] I. Masuda, N. Hagita, T. Akayama, T. Takahashi, and S. Naito, "Approach to a smart document reader system," *Proc. CVPR*, San Francisco, 1985, 550-557.
- [8] D. Niyogi and S. N. Srihari, "A rule-based system for document understanding", *Proc. AAAI*, Phila., 1986.
- [9] K. Kubota, O. Iwaki, and H. Arakawa, "Document understanding system", *Proc. 7th ICPR*, 1984 612-614.
- [10] C-H. Wang and S. N. Srihari "Object recognition in structured and random environments: locating address blocks on mail pieces", *Proc. AAAI*, Phila., 1986.