

DOCUMENT IMAGE UNDERSTANDING

Sargur N. Srihari

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

ABSTRACT

A digital document image is an optically scanned and digitized representation of a printed page that consists of blocks of text, line drawings, half-tone pictures, and icons. Document image understanding is a goal-oriented problem involving detecting and interpreting different blocks and coordinating the interpretations to achieve an end result. We examine several solutions to subproblems in document understanding tasks. The subproblems range from pixel processing issues to symbolic linguistic reasoning to global control.

1. INTRODUCTION

A *document image* is a visual representation of a printed page such as a journal article page, a magazine cover, a newspaper page, etc. Typically, it consists of blocks of text, line drawings, words, and sentences that are interspersed with half-tone pictures, line drawings, and symbolic icons. A digital document image is a two-dimensional array representation of a document image obtained by optically scanning and rasterizing a hard copy document. It may also be an electronic version that was created in that form, say, for a bit-mapped screen or laser printer. An example of a document image is shown in Figure 1.

Document image understanding (or simply, document understanding) is a goal oriented problem that may involve recognizing photographs (which is a *vision* problem), interpreting text (which may involve *natural language understanding*), interpreting line drawings such as graphs, in such a way that the interactions of the different components is documented for. Some examples of document understanding are: (i) integrate the pictorial content of a newspaper photograph with accompanying narrative title, (ii) determine the destination address on a colorful magazine cover: it is necessary to locate and orient the most plausible textual block, to read and verify that it has the correct syntax, and (iii) make inferences from an annotated x-y plot.

Document understanding implies the correct recognition of one of its constituent components. The recognition has to be done in terms of high level components such as columns and paragraphs, as well as in terms of low level components such as words and lines in the case of text; curves, rectangles, and solids in the case of graphics; and image regions in case of half-tone images.

Document image analysis is the task of deriving a high level representation of the contents of a document image. The spatial layout of a document can be represented in several ways. The representations of the image of Figure 1 are shown in Figure 2. The first is an analogical representation that gives the spatial layout of blocks including paragraphs and the second is a semantic representation in the form of a partial semantic network whose nodes represent document entities and arcs represent relationships. Each of them assumes blocks as units of analysis. Each block can be associated with details like size, shape, color, and black pixels and texture information. This

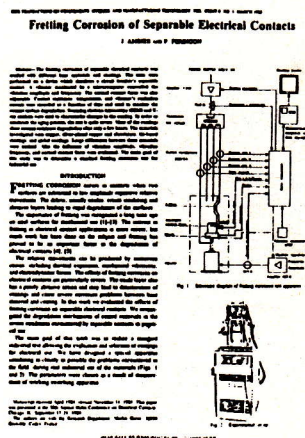


Figure 1. An example of a document image.

information can then be utilized to discriminate between blocks. Methods for deriving the blocks can take advantage of the fact that 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.

Document Understanding Components

Document understanding involves the use of several types of knowledge including *visual*, *spatial*, and *linguistic*. In determining objects from the background, say, by thresholding or by edge detection and grouping, the parameters necessary to perform such operations constitute visual knowledge. Labeling regions involves the use of spatial knowledge, e.g., the layout of a typical document. Determining the font and identity of characters also involves spatial knowledge, i.e., the structure of textual symbols, words, and the distinguishing characteristics of each font. Reading words in blurred or poorly registered text, is a process that involves spatial as well as linguistic knowledge, e.g., lexicon of acceptable words. Determining the role of a block of text, e.g., "is this a title?" is a process requiring spatial, syntactic as well as semantic knowledge. Considerable interaction among different types of knowledge may be necessary. For instance, assigning a role to a textual region may require not only knowledge of spatial layout, but also an analysis of its textual syntax and semantics, and an interpretation of neighboring pictorial and iconic regions.

The various processes necessary for document analysis are shown in Fig. 3. The task has been divided into three phases. Phase 1 consists of two steps: optical scanning and 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. Document understanding involves feedback paths not shown in Fig. 3.

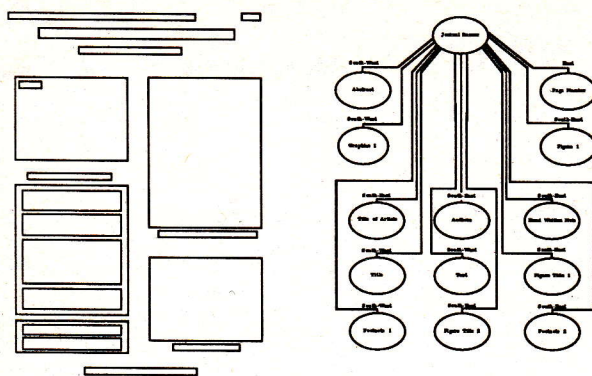


Figure 2. Document image representations.

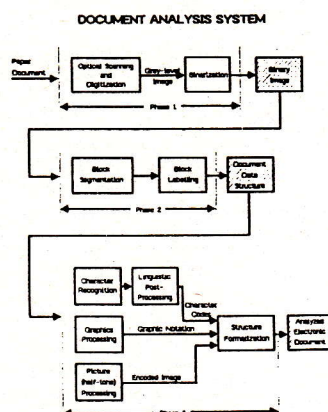


Figure 3. Document Understanding Components.

2. SCANNING AND BINARIZATION

An optically scanned document image is usually an integer (gray-valued) array; it is obtained by a process of spatial sampling and simultaneous conversion of light photons to electric signals. A high-resolution *bit map* is sufficient to capture shades of gray in the eye of the human perceiver (such images are called *half-tone*). An optically scanned document image can be converted into a bit map by a *global* thresholding operation: pixel values below the threshold are deemed to be black (value 1) and those above deemed white (value 0). The threshold itself can either be predetermined or calculated from the image histogram, e.g., the valley of a bimodal histogram.

A global threshold, i.e., one that applies to the entire image, may sometimes be impossible to obtain due to varying background color. One solution is *adaptive* thresholding which is to compare the gray value of the pixel to the average of the gray values in some neighborhood about the pixel. If the pixel is significantly darker than the neighboring pixels, it is called black; otherwise it is white. One such adaptive thresholding algorithm [14] that is appropriate for a printed document considers a 9×9 region around a pixel (x, y) . The central 3×3 region including (x, y) is called A_1 and four 3×3 regions 8-connected to A_1 are jointly called A_2 . The average pixel value of A_1 is compared to the weighted average pixel value of A_2 . The idea in choosing the regions in this manner is that of a black pixel that is part of a line of text should only be compared to the white spaces above and below the line.

Color

A document with significant colored regions is usually

scanned as three gray level images with red (R), green (G), and blue (B) filters. A colored image can be converted to a bit map by approximate methods. Color thresholding to extract regions of a particular color. Thus if the goal were to determine the address label which is known to be red, then one can focus only on the white regions and ignore the further feature tests (e.g., rectangularity, presence of

Resolution

Scan resolution is of importance in document understanding. Scanning resolution for text and line drawings can be decided as follows. The width of a typical stroke is about 0.2mm (0.008 inch), with some of the strokes up to about 1mm; a 10-point character means 0.5mm (0.014 inch) between ascender and descender lines. A sampling rate of 240ppi corresponds to about 0.01mm which guarantees that at least one pixel will fall totally within a stroke; this resolution corresponds to 33 pixels between ascender and descender lines. In the case of line drawings the minimum line width is approximately 0.1mm. Thus the minimum resolution needed is approximately 240ppi.

3. BLOCK SEGMENTATION

Approaches for segmenting document image can be roughly divided into techniques that are top-down and bottom-up. Top-down techniques divide the document into major regions which are further divided into sub-regions. They often use knowledge about document structure. Bottom-up methods progressively refine the data by layered operations. Although no practical system takes a purely bottom-up approach, every system can be identified as being predominant with one of the two philosophies.

3.1. Top-Down Approaches

We will describe three top-down approaches: projection profile cuts, and the Hough transform.

Smearing

The smearing technique [24] involves scanning a document as black/white images. A run-length smearing algorithm then 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 are left unchanged. For example, if the input sequence is 00011000001100100001 and the value of t was 3 then the output of the RLSA on this sequence would be

11111000001111100001.

The RLSA is first applied row-by-row and then column-by-column, yielding two distinct bit maps. The two are then combined by applying a logical AND to each pixel. The resulting RLSA image contains a *smear* where the material appears on the original image. The threshold in the two directions need not be the same. The segments expected to yield blocks each of which should contain a particular type of data (text, graphics, half-tone, etc.).

Projection Profile Cuts

Based on the observation that printed pages are made up of rectangular blocks, a page can be recursively tiled into rectangular blocks. Thus the document is represented as a tree of nested rectangular blocks. The application is based on the configuration of the pixels. A "local" peak is applied to horizontal and vertical "profiles" to detect peaks (corresponding to thick black or white gaps) at which cuts are placed; it is local in that the width is determined at the nesting level of the recursion, e.g., gaps between paragraphs are thicker than those between lines [11, 25].

A complete hierarchical recursive tiling of a page

represented as a tree. An *X-Y tree*, is obtained by using horizontal or vertical partitions at alternating levels of the tree. The root represents the entire page and can be considered either a single horizontal or vertical partition. The succeeding nodes, and finally the leaves are obtained by a process of alternating horizontal and vertical cuts. The number of horizontal or vertical cuts placed at each level is variable, so the resulting tree is not binary. A node that represents a horizontal (vertical) partition has a set of descendents 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. It involves transforming each black pixel (x, y) in the original image into a curve in the parametric space. In the case of detecting straight lines the transformation used is $K = x \cos T + y \sin T$. For each point (x_i, y_i) , T is varied from 0 to J to yield the corresponding curve in the (K, T) space. An array of accumulators is set up by quantizing values of K and T . The accumulators corresponding to values of (K, T) yielded by the transformation are incremented by one. It can easily be shown that if points are collinear along the line specified by parameters (K_0, T_0) , then each such point (x_i, y_i) will necessarily increment accumulator (K_0, T_0) . Thus each value in the accumulator array corresponds to the strength of evidence for a straight line with the corresponding parameters. For a 512×512 image, K values extend from -368 to $+368$ (This value is arrived at by assuming the origin of the (K, T) space to be at the center of the 512×512 image and hence the maximum value of K is 256×2).

The accumulator array resulting from applying the Hough transform to a portion of the document image of Figure 1 is shown in Figure 4. The portion corresponds to the "abstract" in the left column. 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. In the case of text containing both upper and lower case, they register a maximum value at the center line of the characters and slightly lower values corresponding to the ascender and descender lines. Textual lines have another property that there are "white" lines perpendicular to the beginning and ending columns and this information can be used to confirm line orientations. First, the number of vertical columns in a page needs to be determined because a multi-column page may not have any (K, T) orientation with low counts. Subsequently, the spacings between the text lines can be identified, on the basis of the average width of significant (K, T) values for a particular T .

The distance histograms along the line and its perpendicular orientation can be used to get the inter-word and inter-line distances. By incorporating these, the thresholds for the run-length smearing algorithms can be derived and big blocks of data can be grouped together [16].

Bottom-Up Approaches

Bottom-up analysis is based on successive steps of refinement of the input image. The image is first processed to determine the individual connected components. At the lowest level of analysis, there would be individual characters and other line 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. In a bottom-up technique 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.

Size

A commonly used method is to determine the size of a component, by counting the number of pixels, and consider it not to be a character if it lies within a range. A disadvantage of this method is that the size of characters has to be known beforehand, which is not always possible. However, size is a useful first feature.

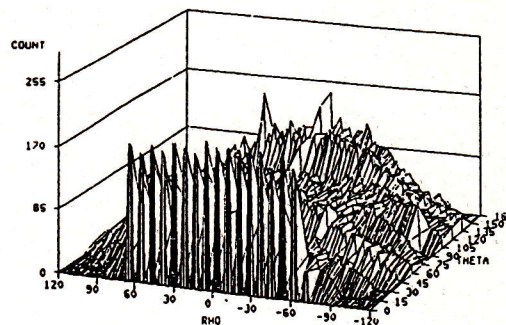


Figure 4. The Hough Transform accumulator array. Text lines appear as a series of peaks in the 0° column.

Branching Structure

The geometrical complexity of an object can be measured in terms of 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, i.e., it contains too many branching structures to be a printed character. An algorithm can be designed to determine geometrical complexity of several components while the image is scanned raster fashion. This is done by keeping track of the several components in each line and their connectivity with respect to components in the previous line. A description of each connected component consisting of, say, its size, location, horizontal and vertical extents is maintained. Clusters of structures of similar type can then be identified. An example of this clustering approach is shown in Figure 5. A pass is made on the binary data to determine all connected components, which are individual characters. On the basis of proximity, size and other features the characters are merged into lines [2].

Topology

The 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). A fast algorithm can be designed for computing the Euler number of a binary image. More generally, the adjacency tree of a binary image contains all topological information (Figure 6).

Control for Bottom-Up Grouping

The order of application of different operators can be based on their complexity, e.g., most characters can be quickly determined by their component size. However if characters touch a line, as often in the case of line drawing graphics, the characters have to be segmented from lines. One technique for segmenting characters from line structures is to determine the high neighborhood line density (NLD) areas in the line structures [8]. Such areas are candidates.

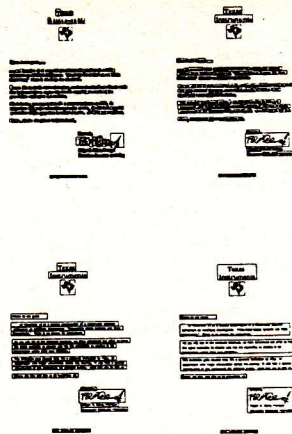


Figure 5. Bottom-up grouping technique. Character blocks are merged into words, words into lines, etc. using a clustering technique.

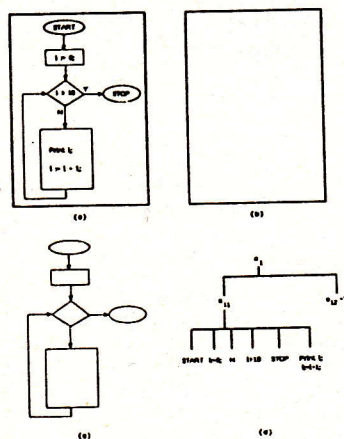


Figure 6. Adjacency tree: (a) input graphics, (b) outermost component a_1 , (c) inner component a_{11} , and (d) complete adjacency tree.

The application of different operators for bottom-up grouping can be coordinated by a production system. Examples of specific rules in such a system are:

- if connected components $[c]$'s size $[s]$ is larger than t_1
- then the component is graphics with likelihood L_1 ,
- if neighborhood line density (NLD) on graphics is high,
- then the high NLD area is a character area with likelihood L_2 .

A sketch of control flow follows:

1. Size, position, and direction data of every connected component is collected.
2. The components are classified according to size into graphics, characters, and noise and their likelihoods are determined.
3. Character components are grouped together as character areas using vector distance rules.
4. High NLD areas in graphics components are separated out as characters touching lines and their likelihoods are determined.
5. If another character exists near the touching character area its likelihood value is increased.
6. The cut-off part in a faded line is repaired.

4. LABELING

Blocks determined by a segmentation process usually need

to be labeled in a subsequent step. There are several techniques for classifying a given block into one of a small set of predetermined document categories.

Statistical Approach

Simultaneously with component coloring, the following measurements are taken: total number of black pixels in segmented block, minimum x - y coordinates of a block and x - y lengths, total number of black pixels in the original image for the block, and number of horizontal white-black transitions in the original image block.

The next step is to classify each block according to content. This is done by computing several features for each block from the above measurements and then using a linear pattern classifier. The features computed are: the height of the block, its eccentricity, the ratio of the number of black pixels in the block to the area of the surrounding rectangle, and the mean horizontal length of the black runs in the original data from the block.

A block is determined to be text if it is a textured shape. The mean height H_m and mean length of black run R_m are used. The distribution of values in the R - H plane derived from sample documents are observed to determine the discriminant function. Low R and H values represent regions containing text. To determine the threshold values of R and H that define the region in the R - H plane, an adaptive method is used. This method estimates R_m , H_m and the standard deviation of R and H . These values are then used to classify the various blocks using the following pattern classification scheme that assumes linear separability:

- **Text:**
if $R < C_{21} \times R_m$ and $H < C_{22} \times H_m$,
- **Horizontal solid black lines:**
if $R > C_{21} \times R_m$ and $H < C_{22} \times H_m$,
- **Graphic and Half-tone images:**
if $E > 1 / C_{23}$ and $H > C_{22} \times H_m$ and
- **Vertical solid black lines:**
if $E < 1 / C_{23}$ and $H > C_{22} \times H_m$.

A Distance Mapping Shape Measure

A shape measure for determining whether a component is a character, line drawing, or thresholded gray-level image follows [22]. Given a connected component S first a distance is computed for each pixel of S as follows. Distance is a function of the two Cartesian coordinates x, y ($(x, y) \in S$) and an angle ϕ measured from the x axis ($0^\circ \leq \phi \leq 180^\circ$). It is defined as the length of a line segment B_1B_2 that passes through point (x, y) of S and makes an angle ϕ with the horizontal. Two opposite border points B_1, B_2 of S are such that the segment B_1B_2 is entirely inside S . Three distance mappings

$D_{min}(x, y) = \min_{\phi} [d(x, y, \phi)]$, or the length of the shortest chord through (x, y) ,

$D_{max}(x, y) = \max_{\phi} [d(x, y, \phi)]$, or the length of the longest chord through (x, y) , and $D_{ecc}(x, y) = D_{max}(x, y) / D_{min}(x, y)$.

These mappings reflect geometrical properties of shapes within binary images. By averaging D_{min} within S the minimum border to border distance d_{min} can be calculated. It provides a fairly good estimate for the mean line thickness of line shaped patterns. To discriminate between line shaped patterns and more compact ones, the average of D_{ecc} is a useful measure. Similarly, d_{max} reflects the mean maximum border distance.

Further shape factors are desirable: d_{min} can be used to derive a shape factor which equals one for the most compact

pattern, an ideal circle, and is greater than one in all other cases: $f_1 = A / d_{min}^2$ with $C_1 = 16 / 9\pi$ in continuous space (or number of pixels in discrete space). This global measure reflects the "line shapeness" of pattern, i.e., their compactness independent of A more specific characterization can be obtained by means of additional features, e.g., $f_2 = C_2 \times A / d_{max}^2$ with $C_2 = 4 / \pi$ in continuous space (1.2.388 in discrete space) which is invariant with respect to the size of the pattern. Unlike f_1 , f_2 is capable of discriminating a straight line from an L-shaped line having about an equal number of pixels and about equal line-thickness.

Handwriting versus Print Discrimination

Given a region that contains either printing or cursive script, a method can be devised to determine its nature by observing regularities of white-to-black and black-to-white transitions. Figure 7 illustrates.

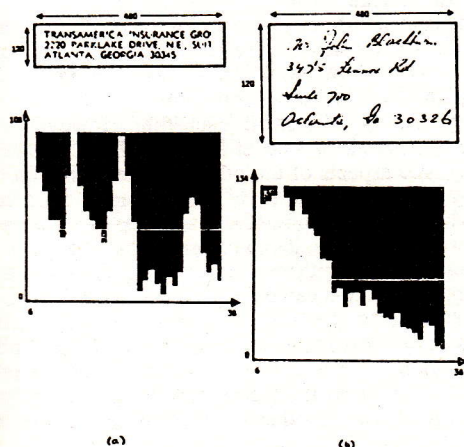


Figure 7. Printed text can be distinguished from cursive script by determining regularities: (a) printed text and histogram of black-white transitions in vertical direction, (b) cursive script and corresponding histogram.

The method 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 . Different histograms are obtained by considering a white-to-black transition separated from a black-to-white transition in the horizontal/vertical directions. The histograms for printed text tend to have several regular peaks while handwriting tends not to have such peaks. Thus the test becomes one of classifying the histogram as regular/irregular. The classification can be done either by treating the histogram itself as a feature vector or by extracting an appropriate feature vector from the histogram.

Rule-Based Approach

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. This can be done by a rule-based expert system that uses a knowledge-base comprising layout and composition rules for specific classes of documents [11]. Using physical (syntax) information, the rule-based system proceeds on its labeling of the blocks (referred to as semantics). These rules can determine line-of-characters, paragraphs, columns, photos etc. using concepts of width, aspect ratio, number of cuts, number of siblings at a particular stage and other such related information. Rules are of the

following type:

- line of characters: sequence of adjacent character blocks of some height, separated by character-segmenting rule,
- paragraph: sequence of blocks of line-of-character of same length, separated by line finding rule,
- column: large blocks with approximately equal frequency of "1" and "0" pixels, and
- photo: sequence of paragraph blocks of same width, separated by paragraph-cutting rule.

For a complete consistent set of rules, the labeling process has worst case complexity equal to the product of the number of blocks and the number of labels.

5. GRAPHICS PROCESSING

Line drawings and tables are commonly encountered in documents. Their analysis involves a raster to vector conversion. The aim of a raster to vector process is to convert a binary pixel representation of line-work into a connected set of segments and nodes. A segment is typically a primitive such as a straight line or a parametric curve. Straight lines in vector form are specified by their start positions, extent, orientation, line width, pattern, etc. There have been several different approaches to vectorization, including: pixel-based thinning, run-length based vectorization, and contour-based axis computation.

Pixel-based thinning algorithms require multiple passes through the data set. During each pass object pixels are deleted based on local neighborhood criteria. Such methods gradually thin thereby thinning the image down to unit width (Figure 8). A single pass algorithm could be obtained by using a set of local 3×3 operators and applying them in a quasi-parallel manner to generate a marked skeleton; quasi-parallel application means dividing the image into four disjoint sub-images and applying the operators to each sub-image in turn [5]. These techniques are slow as they operate on individual pixels rather than on groups of pixels. Also a lot of time is wasted in looking at white pixels more than once; only about ten percent of the pixels on any page of text are black and hence this wastage is quite considerable.

Run-length based vectorization [15] differs from pixel-based thinning in that the image does not go through multiple passes to get at a thinned image of unit width. Instead, vectors are drawn with certain approximations resulting in this. The method aims at this by tracking the black segments along scan lines and using aggregates of such segments to operate upon. The algorithm achieves the objective by converting the input run length encoded image into a graph and analyzing the graph. The graph also has information which could be used to provide many more useful features for character recognition. It is found that this method works better for lines perpendicular rather than parallel to scan lines. The complexity of pixel based thinning is of the order of the number pixels. The complexity of run-length based thinning is approximately of the order of the number of run length code segments.

Some problems that need to be specially examined are handling of dotted lines, determining end points of lines, effectiveness of thinning, determination of vector intersection and curve fitting versus piecewise linear approximation. Issues in understanding of graphics are discussed in [3, 4, 10].

6. TEXT RECOGNITION

A text recognition or reading technique takes an image of text and maps it into an ordinal machine representation, e.g., ASCII. Methods of reading text can be divided into two related but distinct categories: *isolated character recognition* and *character recognition in context*. When the types of font are known *a priori*, or when few fonts are encountered, a highly reliable isolated character recognition technique can be designed.

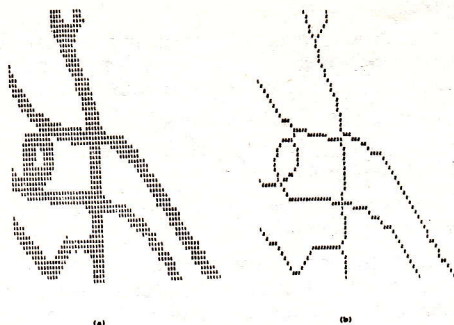


Figure 8. Example of Thinning Line Drawings.

Isolated character recognition involves extracting a set of predetermined features from the character image. The features used are chosen for their ability to discriminate between classes. For instance, character height in millimeters is useful for distinguishing between lower case letters *with* ascenders, (e.g., b, d, h, k, l, t) and letters *without*, (e.g., a, c, e, n, o, r, s, u, v, w, x, y, z). Thus the character is regarded as a point x in a multi-dimensional feature space. Classification is accomplished by using discriminant functions which effectively partition the feature space; typically each class i has an associated discriminant function d_i and x is assigned to the class that maximizes $d_i(x)$. The discriminant functions d_i , or their parameters, are determined statistically from a large set of samples during a training phase. Character recognition techniques quite often use only binary-valued features, since they lead to simpler feature representation and discriminants

6.1. Isolated Character Recognition

Typical of such a method is one [17] that treats a binary character image on a 16×16 grid as a 256-element binary feature vector x . The classifier uses second order polynomial discriminant functions derived using least mean squared error considerations. The form of the discriminant function is: $d_i = a_{0i} + a_{1i}y_1 + a_{2i}y_2 + \dots + a_{mi}y_m$, where y_i are either components of x or products of pairs of components. Since the number of such terms is large, (i.e., $256 + {}^{256}C_2$), only a subset is used. Approximately $m = 1200$ is determined where the pairs of pixels are heuristically chosen. In the 16×16 raster field, more pairs are chosen at the center than in the periphery. The set of coefficients a_{ji} are represented as a matrix A of size $k \times m$. Thus the computation of the k discriminant functions can be expressed as the matrix multiplication: $d[k \times 1] = A[k \times m] \cdot y[m \times 1]$. The index of the member of d which is closest to unity is the class to which the input is assigned.

6.2. Character Recognition in Context

In most printed documents characters hardly ever appear in isolation. Several approaches to recognizing characters in context are known [19]. Contextual information is usually in the form of a lexicon of acceptable words. It can also be in the form of n -grams, i.e., legal letter types, or letter transition probabilities.

One approach is to recognize characters of words and then post process (correct) any errors by using a lexicon of words and a distance measure between words. Another approach is to use lexical information to limit the number of possibilities in recognizing individual characters to find the nearest word. This is done by weighting the choices for a given character by (i) its occurrence in the lexicon and (ii) its frequency of occurrence in the text. Thus the unlikely choices are eliminated and the

resultant is guaranteed to be in the lexicon. Contextual information into the character recognition leads to better performance than the two step character recognition followed by error correction. A third approach is to extract features from the character attempt to classify it using a lexicon organized by a simple set of features are used in a first level and a neighborhood of words and a more detailed discriminates between a small subset of characters.

6.3. Type Font Analysis

A document analysis system needs to distinguish text of different fonts. In the terminology of typography is a collection of upper and lower case letters, characters of one particular typeface, style, and typeface determines the overall design of one character.

The style refers to the average stroke width of boldface versus lightface (normal), and the position italic versus roman.

The task of character recognition is considered if the font is determined ahead of recognition. Characters is due to two sources: the type font and the English alphabet in several fonts, the style of the letter, say "B," is different in each font. However, one type font certain aspects of the style of the letter to the stylistic aspects of the other characters in the lower left sides of "B" and "D." By providing within which this repetition of stylistic aspects can be expressed explicitly, it can be used subsequently in the recognition algorithm, e.g., knowing the procedure for that font can be applied.

Determining the stylistic consistency of characters involves finding elements such as serifs and the relationships to strokes to which they are joined. This is for type font analysis considering only machine taken from the Roman and sans-serif families in

6.4. Linguistic Analysis

The role of language in the recognition of characters has been little explored beyond the word level. To describe a linguistic approach to recognize words, a dimensional block of text is an address. If a parse tree is obtained. The method exploits spatial information in its analysis, i.e., the space between words is important.

In determining the destination address of a character image with several two-dimensional blocks of text necessary to syntactically analyze the recognition candidate block by using spatial information of word and line spacing [21]. A system has been developed to perform this syntax analysis by attempting to subelements of the block of text and producing a representation of the labeling of the constituents [18]. This identifies words corresponding to the city, street name, among others. Several significant changes have been made to a standard Augmented Transition Network parser [18] including word abbreviation analysis, and heuristic parsing.

This module takes as input a sequence of characters (subsequently called a *sentence block*). This block of ACSII characters, is then preprocessed to a sequence of horizontal blocks before each word and line information as lexical features of the input. This preprocessing operation is performed to remove the spacing between lines in the block but with the word recorded as a special word inserted directly into the input and later handled as another word. This provides the parser with spatial information to check if a word or line is "too far away" from

words to be considered.

This preprocessed sentence block will be *accepted* by the parser if and only if the sentence block has the syntax of an address block. If the input sentence block is acceptable in terms of the grammar and lexicon, then the structure (parse tree) of the region will be passed to the generalized rule-based region labeler for use in the determination of the correct sending address as opposed to the return address.

The address recognition module incorporates the four assumptions described below:

1. When a block is initially segmented and passed to this address recognition unit, extraneous printed material may be present in the block. This is dealt with by adding heuristic arcs (described below) to the grammar which will allow crude matching of the beginning of each line and result in the skipping of the remainder of the text line using the generalized TO arc [18].
2. If the text block contains an ambiguous address, then the output will be a parse tree representing all possible interpretations of the block.
3. This unit will not perform context matching between the city, state, and street name with the five or nine digit zip code. Context checking is also not performed between other components of the address since it is not a necessary condition for a text block to be an address block.

Parser Modifications

The major changes to the standard parser include the capturing of spatial features such as end of line and word spacing, and the addition of abbreviation handling.

Spatial Properties: A major parser modification was the addition of a spatial acquisition mode for capturing the number of spaces before each word and adds this spatial information to each lexical form of the current input word. This information includes the number of blank spaces prior to a word, the distance to the left, top, and right margins and also the length of the word. Below is a partial input to the parser after the initial conversion of physical spaces in the input to lexical features associated with the words.

```
(buffalo(((root . buffalo)
(city . city)
(parts . 1)
(head-space 3)
(left-margin-space 3)
(top-margin space 0)
(word-length 7)
(right-margin-space 29)
(line-length 32)))
```

If vertical spacing is present in the input, then the block will be identified with the insertion of a new line *word* of the form: *@@@* where @ represent the start of a new line, and the *nnnn* represents the number of units (lines) in the vertical spacing.

Abbreviations: Since many of the words in an address block are abbreviations such as Ave, St, NY, PA, S, N, W, a method was developed to associate these abbreviations with their expanded forms. To facilitate the lexical look-up of abbreviations and avoid the repetition of lexical words with similar meaning, a new property was used to equate a word with names of other words.

A grammar corresponding to the bottom two lines of a standard address has includes the structures of the PO box, street, city, and state fields as well as the five (or nine) digit zip code. This grammar is complicated by the possible presence of punctuation between words and the spatial information mentioned above. Sometimes the punctuation and spatial information is important. For instance, if a comma is present,

then it is a separator between the city and state name. Also, if a new line word is present, then it is most likely the termination of a multiple word phrase such as a city, state, or street name.

The grammar, in addition to containing arcs representing the bottom two lines of a standard address also contains heuristic rules for detecting and skipping other possible input lines. That is, if the first word on a line is a personal title (Dr., Mr., Ms., Miss, or Mrs.), or a personal first name, then the remainder of the line would be ignored. Another skipping rule is for the detection of a standard presort codes usually present on the top line of an address block.

For a sentence block to be parsed as an address region, several constituents must be identified. For instance, if more than one line is present in the address region, then the city/state/sip line must start on a new line. Another necessary element for a block to be parsed as an address block is the presence of a state name followed by a zip code (either on the same input line of on the following input line).

The grammar, initially, makes one pass through the input, attempting to locate a possible state name. If a state name is not found, then the parse will fail. If, however, a state name is found, then the state name and line number are recorded for later processing and a second pass is made through the input utilizing a more robust grammar for detection of other constituents. This preliminary detection of the state name provides an efficient means for detecting a necessary component of the address and reduces the computational cost of parsing the entire address region.

Lexicon

The lexicon contains generally five different types of entries corresponding to the street, city, and state names, street keywords (such as avenue, street, circle, parkway), and other keywords most often present in the address region of mail pieces such as personal names and titles. As was mentioned previously, the lexicon need not contain all possible street and city names but parses of an input sentence blocks containing these "missing" cities of street names would result in weak parses.

The grammar entails 106 states, 178 arcs, and a small text lexicon which contains 188 unique words with entries for all 50 states (and their abbreviations), 32 cities, 20 street names, and 20 street keywords. To judge the effectiveness of the keyword approach in guessing a city name, the Directory of Post Offices (DOPO) data base was used. A total of 25,630 on word city names were encountered with over thirty percent (7,843) of these city names ended in the above mentioned letters. For two word city names. 7,437 cities were encountered with over seventy-four percent (5,522) of these having first or last words as those mentioned above for two word city names. This shows that these patterns are a good indication of a city name, however, testing was also done to compare the number of English words (advertising material) that would be wrongly mapped into a city name based on these rules. Using the Brown Corpus containing 43,264 English words, a total of only 79 (0.18%) English words matched the one word patterns. A comparison of the two word pattern was not possible with the lack of contextual information in the corpus.

The parses of an input sentence which has been accepted using the grammar and lexicon is shown below, where the preprocessed input sentences are also stored in the parse tree as the first element using the form mentioned above.

```
(address(input-sentence
(1241| main @2 n |.| java n |.| y |.152102| |.4128|)
(line starts (00)))
(spatial (word (0 W4 1 W4 0 W1 0 W1 1 W4 1 W1 0 W1
0 W1 0 W1 1 W5 0 W1 0 W4))
(line (L9 2 L23))))
(prepass (states((state (state-name (new york))
```



```

      (puncts (|.| |.)))
(level 01)))
((street-line (street-number |1241|) (direction nil)
  (street-name main) (street-keyword nil)
  (apartment nil) (puncts nil) (street-guess nil))
(bottom-line (city (city name north java))
  (parts 2 (puncts (|.| nil)) (guess t))
  (state (state-name new york)) (puncts (|.| |.)))
  (zip (zip-location after-state)
    (ninezip |151024128|) (fivezip nil)))).

```

7. PICTURE PROCESSING

A document image understanding system needs to integrate the information contained in photographs with the information contained in other fields such as with the text and line drawings. Often an understanding of the photographs is critical to understanding the narrative. Deriving descriptions of the contents of photographs is the subject of many computer vision projects today.

8. CONTROL STRUCTURES

The coordination of several processes is of central importance in a document understanding system. The control structure allows the application of the appropriate process. For instance, there may be more than one way of binarization. The method to be applied is determined by the control structure.

Several efforts towards designing control structures for document understanding are described [8, 9, 12, 23]. A method based on production systems and blackboard communication first proposed in the context of speech understanding holds promise for document understanding.

The mechanism described in [8] is based on a production system. When the production system needs data the mechanism causes the appropriate module to start collecting data. In this mechanism each connected component is treated as a data unit.

An array, which is a specialized *blackboard* is used for the read/write area in both the production system and program modules. The blackboard is made up from three planes of identically sized arrays. Data, indicators, and program modules names and arguments are respectively stored in these planes.

Data concerning positions, size, likelihood, and attributes for connect components are stored in the first data plane. The indicator plane indicates whether or not the value in the corresponding position in the first plane has already been collected. When the production system requires some data for matching, the indicator is checked first. If the desired data has already been collected, the values in the data plane are transferred to the production system. In the case where the data has not yet been collected, the program module stored in the third plane is activated and its function is carried out. The program module function determines the data and forwards the values to the data array.

Advantages of using production systems in document understanding are:

- easy to apply either an additional technique to an object that is hard to interpret with only one technique, or a retry process having modified parameters,
- processing knowledge is expressed as rules, thus software maintenance (e.g., modifying or patching of a program) becomes easier, and
- as the processing is not carried out over the entire picture uniformly, but only in necessary segments, high efficiency is obtained.

An Expert System

A production system that is organized with different levels of production rules that perform an analysis of a document image, and interpret and classify the various regions of printed matter on the document is given in [12]. The control flow is as

strategy is to be followed at any given time for analyzing the image. This means that the strategy rules determine what the order of execution of the control rules will be.

Representation of Uncertainty

The system has to deal with many situations where a combination of rules, rather than a single rule, lends credence to a particular hypothesis. Thus, the success of each of these rules is measured in terms of evidence towards that hypothesis. If the total evidence obtained from the successful rules is sufficiently high, then the hypothesis is assumed to be true, and the next stage in the process can then be tackled with the assumption that the hypothesis has been confirmed. To deal with such a situation, each Knowledge Rule in the system is given a certain confidence value between 0 and 1. When the knowledge rules are applied, the confidence values for all the rules that succeed are summed. The sum thus obtained indicates the *certainty factor* of the conclusion obtained from the control rule which invoked the knowledge rules. This certainty factor is used for the purpose of ordering the conclusions at any given stage so that the more likely conclusions can be examined in further detail than the less likely ones. This has the effect of making the process more efficient, thus reducing execution time in many cases.

Conclusion

The wisdom of using production rules to represent knowledge has been the topic of discussion among many AI researchers. In the domain document understanding, a rule-based system is extremely elegant because unlike natural scenes, documents are very structured in character, and thus knowledge about the features of documents can be very effectively formulated in terms of production rules. There are other advantages to using production rules in document image understanding. First, it is easy to apply either an additional strategy to a region that is difficult to interpret with only one strategy, or a retry process using modified parameters. Second, software maintenance is made easier, since addition/modification of rules is a relatively simple process that does not disrupt the rest of the system. Third, in a production system the processing is not spread out over the entire image uniformly, but only on relevant segments; thus, high efficiency is achieved. All these advantages make production rule-based systems eminently suitable for use in the domain of document understanding.

DISCUSSION

Document image understanding is a task that begins with image processing and ends with complex symbolic reasoning. It is an area of research that draws upon techniques of image processing, pattern recognition, computer graphics, natural language processing, and artificial intelligence. There is a high level of activity in this field in Japan and in Europe. Research in this topic is also growing in the United States.

Several of the components for building a document understanding system are now well-understood, e.g., image detectors, line detectors, single font character recognizers, text parsers, etc. Several other components need to be defined, e.g., techniques for text region determination. Efficient character recognition without operator input will continue to be a challenging problem for the foreseeable future. The problem of raster to vector conversion of line drawings is formidable but several problems remain, e.g., handling of overlapping lines, global considerations in thinning, etc.

The coordination of component processes is a problem that has been addressed in other domains such as speech understanding. The coordination of document understanding processes will have to be done using similar techniques with domain knowledge. Due to advances in several related areas, it can be concluded that document analysis is now a task that is well defined and of a moderate level of complexity. Thus

the prospects of tangible results are reasonably good.

ACKNOWLEDGEMENT

Several members of the Document Image Understanding Group at Buffalo contributed to this work. In particular, Jon Hull, Debashish Niyogi, Paul Palumbo, Joseph Piazza, Steve Tylock, and Ching-Huei Wang. Gregory Zack of Xerox Webster Research Center provided the impetus for this work and many suggestions. This work was supported in part by the United States Postal Service contracts 104230-84-D0962 and 104230-85-M3349.

10. BIBLIOGRAPHY

- [1] C. Cox, B. Blesser and M. Eden, The application of type font analysis to automatic character recognition, *Proc. of Second IJCPR*, Copenhagen, 1974, 226-232.
- [2] W. Doster, Different states of a documents content on the way from the Gutenbergian world to the electronic world, *Proc. Seventh ICPR*, Montreal, 2, 1984, 872-874.
- [3] M. Ejiri, T. Miyatake, S. Kakumoto, S. Shimada, and H. Matsushima, Automatic recognition of design drawings and maps, *Proc. Seventh ICPR*, Montreal, 2, 1984, 1296-1305.
- [4] R.P. Futrelle, A framework for understanding graphics in technical documents, *Proc. IEEE-CS Expert Systems in Government Symposium*, McLean, VA, 1985, 386-390.
- [5] J.F. Harris, J. Kittler, B. Llewellyn, and G. Preston, A modular system for interpreting binary pixel representations of line structured data, in J. Kittler, K.S. Fu, and L.F. Paul (eds.), *Pattern Recognition Theory and Applications*, D. Reidel, 1982, 311-351.
- [6] J.J. Hull and S.N. Srihari, A computational approach to visual word recognition: hypothesis generation and testing, *Proc. IEEE-CS Conference on CVPR*, Miami Beach, 1986, 156-161.
- [7] K. Inanaga, T. Kato, T. Hiroshima, and T. Sakai, MACSYM: A hierarchical parallel image processing system for event-driven pattern understanding of documents, *Pattern Recognition*, 17(1), 1984, 85-108.
- [8] K. Kubota, O. Iwaki, and H. Arakawa, Document understanding system, *Proc. Seventh ICPR*, Montreal 1, 1984, 612-614.
- [9] I. Masuda, N. Hagita, T. Akiyama, T. Takahashi, and S. Naito, Approach to a smart document reader system, *Proc. IEEE-CS Conference on CVPR*, San Francisco, 1985, 550-557.
- [10] F.S. Montalvo, Diagram understanding: the intersection of computer vision and graphics, Mass. Inst. of Technology, AI Memo 873, 1985.
- [11] G. Nagy, S.C. Seth, and S.D. Stoddard, Document analysis with an expert system, *Proc. Pattern Recognition in Practice II*, Amsterdam, June 19-21, 1985.
- [12] D. Niyogi and S.N. Srihari, A rule-based system for document understanding, *Proc. AAAI-86: Fifth National Artificial Intelligence Conference*, Philadelphia, 1986.
- [13] P. Palumbo and S.N. Srihari, Text parsing using spatial information for recognizing addresses in mail pieces, *Proc. Eighth ICPR*, Paris, 1986.
- [14] P. Palumbo, P. Swaminathan and S.N. Srihari, Document Image Binarization: comparison of techniques, *Proc. SPIE Symposium on Digital Image Processing*, San Diego, 1986.
- [15] T. Pavlidis, A hybrid vectorization algorithm, *Proc. Seventh ICPR*, Montreal, 1, 1984, 490-492.

- [16] A. Rastogi and S.N. Srihari, Recognizing textual blocks in document images using the Hough transform, TR 86-01, Dept. of CS, SUNY at Buffalo, 1986.
- [17] J. Schurmann, A multifont word recognition system for postal address reading, *IEEE Trans. Computers*, C-27, 8, 1978, 721-732.
- [18] S.C. Shapiro, Generalized augmented transition network grammars for generation from semantic networks, *American J. of Computational Linguistics*, 8(1), 1982, 12-25.
- [19] S.N. Srihari, *Computer Text Recognition and Error Correction*, IEEE Computer Society Press, Silver Spring, MD, 1984.
- [20] S.N. Srihari, J.J. Hull, and R. Choudhari, Integrating diverse knowledge sources in text recognition, *ACM Transactions on Office Information Systems*, 1(1), 1983, 68- 87.
- [21] S.N. Srihari, J.J. Hull, P.W. Palumbo, D. Niyogi and C-H Wang, *Address Recognition Techniques in Mail Sorting: Research Directions*, TR85-09, Dept. of CS, SUNY at Buffalo, August 1985.
- [22] F.M. Wahl, A new distance mapping and its use for shape measurement on binary patterns, *Computer Vision, Graphics and Image Processing*, 23, 1983, 218-226.
- [23] C-H. Wang and S.N. Srihari, Object recognition in structured and random environments: locating address blocks on mail pieces, *Proc. AAAI-86: Fifth National Artificial Intelligence Conference*, Philadelphia, 1986.
- [24] K.Y. Wong, R.G. Casey and F.M. Wahl, Document analysis system, *IBM Journal of Research and Development*, 26(6), November 1982, 647-656.
- [25] H. Zen and S. Ozawa, Extraction of the fair document from mixed mode manuscript, *Proc. Conference CVPR*, San Francisco, 1985, 544-549.