

The use of document structure analysis to retrieve information from documents in digital libraries

Debashish Niyogi and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo, Buffalo, NY 14228-2567
{niyogi,srihari}@cedar.buffalo.edu

ABSTRACT

This paper describes an approach to retrieving information from document images stored in a digital library by means of knowledge-based layout analysis and logical structure derivation techniques. Queries on document image content are categorized in terms of the type of information that is desired (e.g., articles on a given topic), and are parsed to determine the type of document from which information is desired, the syntactic level of the information desired, and the level of analysis required to extract the information. Using these clauses in the query, a set of salient documents are retrieved, layout analysis and logical structure derivation are performed on the retrieved documents (using DeLoS, a document logical structure derivation system developed at CEDAR), and the documents are then analyzed in detail to extract the relevant logical components. A “document browser” application, being developed based on this approach, allows an user to interactively specify queries on the documents in the digital library using a graphical user interface, provides feedback about the candidate documents at each stage of the retrieval process, and allows refinements of the query based on the intermediate results of the search. Results of a query are displayed either as an image or as formatted text.

Keywords: Document image understanding, logical structure analysis, layout analysis, information retrieval, digital libraries.

1. INTRODUCTION

In document image understanding, syntactic and semantic interpretations of the structure and contents of document images can be obtained using layout analysis and logical structure analysis. This paper describes an approach to retrieving information from document images stored in a digital library by means of knowledge-based document layout analysis and logical structure derivation techniques.

Digital document libraries have become an increasingly important means of storing critical information within organizations. The need for accurate and up-to-date information has led to the practice of storing important information in digital documents that can be made accessible to a wide variety of potential users. This is particularly useful when the information needs to be efficiently represented, updated, and instantly retrievable on demand.

A digital document library typically contains a set of document images that have been indexed according to document type. Additional indexing on documents can be achieved by performing OCR on each document and extracting textual information, but maintaining such a detailed index is often neither efficient nor desirable. One alternative solution is to create a detailed index for a document only when required by a user query, so that only relevant documents are indexed and retrieved. Also, typically a query on a document library yields complete documents that satisfy the query criteria; it is often desirable to extract only relevant parts of a document, a process that is much harder and requires more detailed knowledge of the document structure.

Such detailed knowledge of document structure can be obtained by performing layout analysis and logical structure derivation on the document. Document layout analysis and logical structure derivation enable us to determine the relationship between the *physical* layout of a document page (consisting of the geometric structure and spatial relationships of the different blocks of printed matter) and its *logical* layout (consisting of the logical groupings of related blocks into composite units). This can then be used to actually “translate” a physical document into its logical symbolic representation, and thereby enable the extraction of information contained in a specified part of a document.

Document layout analysis and logical structure derivation (together called document structure analysis) are thus crucial steps in the development of digital document libraries. Extraction of the logical structure of a document enables individual logical components to be identified and indexed separately in digital libraries, thus making their access much faster and easier.

The following sections describe an approach that uses document structure analysis to identify and retrieve information from documents in digital libraries.

2. BACKGROUND

Some ideas that we have developed with respect to the retrieval of data from raster images of documents stored in a digital library have been previously described by Srihari *et al.*¹¹ Related work in the analysis and retrieval of information from document images has been ongoing at CEDAR for several years. For example, retrieval of the identities of human faces in a photograph with the help of its caption has been described by R. Srihari.¹⁰ The determination of the logical structure of a document (i.e., labeling all identifiable blocks, grouping them into logical units, and determining the reading order of the text blocks within a given unit), has been described by Niyogi & Srihari (⁵⁸). Multi-domain document layout understanding has been described by Lam & Srihari.¹ An integrated approach to document decomposition and structural analysis has been described by Niyogi & Srihari.⁶ Techniques for analyzing printed forms, which are examples of complex documents, are described by Niyogi *et al.*⁹ Each of these techniques is being used as part of our approach to analyze and recognize selected components of a raster document image and present the information to a user.

3. DOCUMENT CONTENT RETRIEVAL

We have developed an approach by which document content is retrieved from a set of document images by means of progressively refined levels of document structure analysis. Documents are indexed by performing layout analysis (to extract the syntactic structure of the document) and then logical structure derivation (to extract the semantic structure of the document). The above process yields a hierarchical index of the document structure. This indexing is done at a level that identifies logical units, their spatial relationships, and their components.

4. LOGICAL STRUCTURE ANALYSIS

At CEDAR we have developed and implemented a knowledge-based system for document decomposition and structural analysis. This system, called DeLoS (“DERivation of LOGical Structure”) takes as input the digitized image of a document page and produces as output a symbolic description of the logical structure of the page, including a labeling of each of the printed blocks on the page, a grouping of these blocks into logical units, and the reading order of text blocks within each of these units. Various image processing operations are performed to extract basic properties of the components of the image, and this data is analyzed under the control of a rule-based system, which operates in conjunction with a global data structure to monitor the entire classification, grouping, and block-ordering process. Some details of the design and development of this system have been described in previously published papers (⁷⁵).

4.1. System Organization

Figure 1 shows the components of the DeLoS system. The unbroken arrows between the boxes indicate the flow of data between the various components of the system, and the dotted arrows represent the control flow in the system.

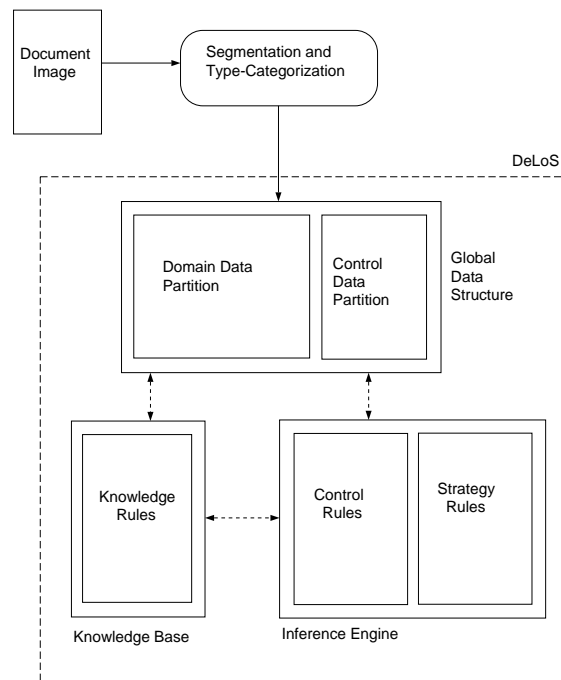


Figure 1: The DeLoS system.

The DeLoS system consists of a multi-level, rule-based reasoning system, an image processing sub-system, and a partitioned global data structure. The rule-based system utilizes a top-down, backward-chaining structure. An inference engine within the rule-based system makes deductions about the document using a hierarchical knowledge base that contains rules describing all the identifiable characteristics of document images. The rules are classified into three levels: knowledge rules, control rules, and strategy rules. The global data structure

facilitates the transfer of information between the image processing modules and the rule-based system. A common data area stores all intermediate computation results and other control information.

The image-processing modules directly access the document image to extract various kinds of information about the document. Intrinsic properties of the different printed blocks as well as the spatial relationships between the different blocks constitute the information that is passed back to the control structure through the global data structure. Some image-processing modules perform the basic tasks of binarization, connected-component analysis, etc., so as to prepare the image for other modules (e.g., block segmentation) that can extract information from it.

4.2. Rule-Based Control Structure

The rule-based system consists of three levels of rules. Using a hierarchical structure of three progressively abstract levels of rules provides a large amount of flexibility in the inference mechanism, and allowed a modular formulation of the solution within the image analysis problem domain. The three levels into which the rules in this system are classified are: Knowledge Rules (level 1), Control Rules (level 2), and Strategy Rules (level 3). The domain knowledge in the system is maintained in the Knowledge rules, and the control decisions in the system are made by the Control and Strategy rules.

4.2.1. The Knowledge Base

The Knowledge rules comprise the knowledge base that contains all the domain knowledge for the system, expressed in terms of first-order predicates. These rules define the general characteristics expected of the usual components of a document image and the usual relationships between such components in the image. Thus, all common characteristics of different types of document blocks (e.g., text blocks, photographs, etc.), as well as spatial constraints commonly followed in document layout (e.g., the positioning of captions relative to photographs, etc.), are encoded into the knowledge base. These knowledge rules in the knowledge base contain all the properties and spatial relations for different types of document blocks, and can be used for block classification, block grouping, or text block ordering as and when required according to the control strategy. Knowledge rules can be further categorized as *unary rules*, *simple binary rules*, and *Complex binary rules*. The hierarchical structure of the knowledge base mentioned above is thus created by these different categories of rules.

4.2.2. The Inference Engine

The control structure for the rule-based system contains an inference engine which is also rule-based, and contains two levels of rules: *control rules* and *strategy rules*. These rules regulate the analysis of the document image, and decide when a consistent interpretation of the image has been obtained. The rules comprising the inference engine are also formulated in terms of first-order predicates. The control structure determines the order in which these rules are executed in order to test various conditions effectively. Control rules regulate the invocation of the knowledge rules, based on appropriate data configurations or processing states. Control rules can be further categorized as *focus-of-attention rules* and *meta-rules*. Strategy rules guide the search in a more general way, i.e., they determine what control strategy is to be followed at any given time for analyzing the image. This means that the strategy rules regulate the invocation of, and determine the execution order of, the control rules. Strategy rules also decide on the stopping criteria for the system, i.e., whether a consistent interpretation and grouping of the blocks in the document image has been achieved (as determined by the absence of incomplete block or unit data in the global data structure, and by the completeness of the logical structure tree). Therefore, there is a set of strategy rules for block classification, another set for block grouping, and yet another for text block ordering.

4.3. The Global Data Structure

The global data structure stores the physical structure and logical structure information for the document being processed. It also facilitates the transfer of information between the rule-based system and the image processing modules. A common data area stores all intermediate computation results and other control information, and provides the framework for the construction of the trees representing the document structures. In this system, the global data structure is divided into the *domain data partition* and the *control data partition*.

5. QUERY PROCESSING

Queries on document image content are categorized in terms of the type of information that is desired (e.g., articles on a given topic, the citation list in a journal article, a photograph with a given caption, etc.). A specific query is first parsed to determine the following clauses:

1. the type of document from which information is desired (e.g., journal page, newspaper page, form, etc.),
2. the syntactic level of the information desired (e.g., a complete article, a text block, a photograph/diagram, a list, etc.),
3. the level of analysis that is required to extract the information (e.g., reading a title/headline, reading the entire text of an article, extracting the contents of a form, etc.).

Using the first clause in the query (i.e., the type of document), the standard indexes for the various documents in the document library are searched, and a set of salient documents are retrieved in the first stage. For example, a query that refers to a journal article will result in the retrieval of only the journal document images from the library.

Using the second clause of the query (i.e., the type of document subset to be retrieved), layout analysis and logical structure derivation are performed on the retrieved set of documents. All documents that do not contain the relevant component are then eliminated from the set (e.g., if the query refers to a photograph, then all documents not containing photographs are eliminated).

Finally, using the third clause in the query (i.e., the level of analysis required), the documents in the retrieved set are then analyzed in more detail to extract the text in the titles, captions, etc., and text understanding is performed if necessary on the extracted text to generate a keyword index. Using this information, the query is refined and the retrieved set of documents searched so as to retrieve the relevant logical unit(s).

6. DIGITAL LIBRARY ARCHITECTURE

At CEDAR, we have integrated multiple information sources to build a digital library of research information on document analysis and recognition. For our digital library, paper documents are scanned for conversion to a form that enables text-based search. The digital library also includes project specific material such as slide presentations, source code, data and video.

Taxila, the CEDAR Digital Library, It collects the results of research and development on various aspects of document analysis and recognition into a unified form for access by all CEDAR researchers. This digital library contains information in several modes, including document images, text files, executable object code, images,

videos, and audio files. The information is organized so as to be retrievable through a World Wide Web interface on a server that is accessible only within CEDAR.

6.1. Organization of Taxila

The CEDAR digital library consists of a research information repository, a set of information sources that classify each piece of information for retrieval, and an information retrieval mechanism through which information is searched for and presented in response to user queries.

The CEDAR digital library consists of the following information sources, described as modules:

- **Image databases** used by CEDAR researchers
- **Interactive computer demonstrations** of systems developed at CEDAR
- **Directory Archives** of information about projects at CEDAR
- **Scanned versions of papers** written by CEDAR researchers
- **Slide shows** of presentations given by CEDAR staff members
- **Source code** developed for various projects at CEDAR
- **Descriptions of techniques** / methodologies developed at CEDAR
- **Tools** that operate on the CEDAR image databases
- **Videos** related to CEDAR research

Each of these modules represent an information source that contains data in multiple modalities. The individual modules are described in more detail below.

6.1.1. Techniques

Innovative techniques and methodologies developed at CEDAR for the solution of various problems in document analysis and recognition constitute a very important part of the digital library. Major areas of research within CEDAR, such as Word Recognition, Japanese Character Recognition, etc., are maintained in this module of the library. For each area, detailed information is maintained, including descriptions of techniques, demonstrations, and reference sources. All of this information is accessible by topic/area. The Techniques module is composed of several sub-modules, corresponding to the major categories into which CEDAR research can be divided. Each sub-module is further divided into specific research topics. Some research techniques are indexed under various sub-modules because of their applicability to different research areas.

6.1.2. Databases

Various document image databases used in ongoing research at CEDAR are maintained within the digital library. These images are in many formats, including Sun Raster, TIFF, GIF, PostScript, HIPS,³ etc. The databases included in the digital library are standard databases of images that have been generated for specific projects at CEDAR. These include databases of handwritten and printed characters & words, images of journal and newspaper pages, etc., and are indexed in terms of their type, source and resolution. The images that are in formats supported by common WWW browsers (*viz.*, TIFF, GIF, and PostScript) can be directly displayed; some images are inlined within the browser window (*e.g.*, GIF) while the browser generates a new window to display others (*e.g.*, PostScript). Images in other formats have to be displayed outside of the browser with the appropriate tools (*e.g.*, the “chips” utility to display/manipulate HIPS format images).

6.1.3. Demonstrations

Demonstrations of systems developed at CEDAR have been incorporated into the digital library. These demonstrations illustrate various document analysis and character recognition techniques that have been developed by CEDAR researchers. Demonstrations in the digital library can be static or dynamic. Static demonstrations show pre-processed results from “canned” demonstrations; this is particularly useful when the individual steps involved in the demonstration are time-intensive, as is the case with many complex image processing operations. Dynamic demonstrations execute a sequence of programs based on inputs given by a user; such inputs are generated from a form that is displayed to the user through a HTML page and filled in and “submitted” by the user to the server which then executes the appropriate programs and displays the results to the user on another HTML page after suitable format conversions.

6.1.4. Directory Archives

Directory archives of information sources for current and past projects at CEDAR are easily retrievable through the digital library. The directory archives are indexed by topic, and each directory archive contains a variety of files in different formats and modalities. For example, a given directory archive may contain a set of program source codes, the corresponding executable object codes, some test images, a text description of the technique used in each of the programs, etc.

6.1.5. Scanned Papers

Scanned versions of research papers are maintained in the CEDAR digital library. The papers are scanned on a Apunix scanner at a resolution of 300 ppi, and the Sun raster files are then converted to GIF format. Users have the capability to display a paper sequentially by page, or to access a specific page of the paper. Also, papers can be printed out on demand for hardcopy circulation if needed. Scanned papers are indexed by topic, and more sophisticated indexing by keywords and concepts is being developed. In addition, other information about each paper is maintained, such as the bibliographic information, including the name of the publication in which the paper appeared, list of authors, date of publication, page numbers, etc.

6.1.6. Slide Shows

Slide shows of presentations are maintained and catalogued in the digital library. Each slide show contains a set of linked slides which constitute a complete presentation. Users of the digital library can either access the slide shows in a sequential manner, or select specific slides from within each slide show and “compose” their own slide show from the selected slides. The latter is particularly useful for researchers who give overview presentations of ongoing research that includes descriptions of techniques developed by various project group members.

6.1.7. Source Code

Source code developed for various projects at CEDAR is an intrinsic part of the digital library. Our objective has been to collect source code from various programs that have been developed to solve certain basic problems, and present them to the researcher who can then modify the code according to the project needs, or simply read the code to gain insights into modular coding techniques. We have collected source code for image processing operations such as connected component analysis, image rotation, image segmentation, layout analysis, etc. The

authorship of each piece of code is specified, and *since the use of the code is strictly internal within CEDAR, no code copyrights are violated*. The code is made available by the author after extensive testing, and is thus only added to the digital library when it has performed the specified task on a large sample data set. Source code included in the digital library can be changed only by the author, and this change is transparent to the digital library retrieval system as long as the file names remain the same.

6.1.8. Tools

Executable tools that operate on the CEDAR image databases are an integral component of the digital library. Many commonly used image processing, vision, and statistical toolboxes and packages already exist on CEDAR machines, and our objective is to make these tools available to all CEDAR researchers through the digital library. By providing these tools to the researchers, we make it possible for them to conduct research in an efficient and productive manner. *System tools* are tools that enable us to perform basic system functions, such as file viewers, compilers, editors, debuggers, etc. — these are generally available. The digital library contains *application tools*, i.e., tools that perform specific operations, such as image processing, vision, neural nets, matrix & vector operations packages, statistical packages, etc., as well as tools developed by CEDAR researchers for solving specific problems in document analysis and recognition. Application tools provide modularity and versatility to the digital library since they execute on data in the document archives. This also results in a “bootstrapping” effect, since the application tools that manipulate the data, e.g., for document analysis and optical character recognition, are derived from contents of the data archives.

6.1.9. Videos

Videos related to CEDAR research are included in the digital library. The objective is to make available multimedia information about different research techniques, as well as their impact on the state-of-the-art in current research and the development of products based on the research. Multimedia presentations about CEDAR research from such diverse sources as scientific television programs and animated illustrations of techniques are stored in the digital library in MPEG or QuickTime format, and when accessed are acted on by the appropriate viewing software that “plays” the video on the user’s workstation. We are also investigating techniques to “index” video sequences, so that the user can go to a specific part of a video based on topic-based annotations created for sequences of interest.

6.2. System Implementation

The server for the CEDAR digital library is a dedicated Sun SPARCstation 2 running the NCSA HTTP daemon version 1.5 under Solaris 5.5, and is connected to a local CEDAR research network of over 150 Sun workstations. Because of this local network connectivity, all information that is stored in the CEDAR shared file systems is directly accessible through the digital library. Since the information contained in this digital library includes results from ongoing sponsored research, the digital library is built as an “intranet” and is not accessible outside of CEDAR.

New information is incorporated into the digital library by first determining the basic information source, i.e., to which section of the library it belongs. This could be considered equivalent to “cataloguing” a new book that is received into a physical library. A primary link is made to the new information through the chosen section, and secondary links are created from other sections that are relevant. For example, when information related to Japanese Optical Character Recognition (JOCR) was added to the library, the primary link was to the Techniques section, where a page was created which gave a brief description of the project, and had links to a JOCR demonstration as well as a link to a description of the CEDAR JOCR database. Secondary links from the

Demonstration and Database sections of the library were then established.

A system has been set up within CEDAR whereby every time a research project, or a logically distinct piece of research, is completed, the researchers involved in that effort use a standard set of specified procedures to convert their research material into the appropriate formats required for the different sections of the digital library, into which the material is then included. Similarly, ongoing research is included in the digital library when it is deemed to be of interest to a significant number of researchers spanning project groups. Each researcher performs the necessary format conversions according to the specified procedures. The researcher is marked as the author of the material when it is converted for HTTP access.

The actual “publishing” of the material, i.e., the incorporation of converted material into the digital library, is done by the CEDAR Digital Librarian, or “cybrarian”, who determines the sections/sub-sections into which each homogeneous unit of information should be placed. This is done in consultation with the appropriate Project Manager who supervised the research. Once the information is placed in the digital library, extensive cross-indexing is performed from every relevant section and sub-section, so that the material is directly accessible from each relevant portion of the digital library.

6.3. Representation of Information

As previously mentioned, each section in the CEDAR digital library represents a specific information source, and there is extensive cross-referencing between the different sections. Each information source is represented by its own “home page”, which contains a brief description of the information source as well as links to all the information accessible in the library that pertains to that information source. For example, the home page for Source Codes contains a brief description of the types of source codes available through the library, and contains links (with descriptions) to individually complete source codes that perform functions such as journal image segmentation, connected component analysis, holistic off-line handwritten word recognition, etc.

Currently, each piece of information is a file in its native format. For example, images are stored as GIF, TIFF, or HIPS format files, research papers are stored as HTML documents or PostScript files, etc. The full range of data formats supported in the digital library is shown in Table 1.

Type of document	Supported formats
Description of techniques	HTML, ASCII text
Scanned papers	GIF, TIFF, HIPS, PostScript, Sun Raster
Source code	ASCII text
Videos	Quicktime, MPEG
Demonstrations	Perl, C, CGI Forms, Java, JavaScript
Directory archives	FTP-able Solaris file structures
Databases	GIF, TIFF, HIPS, PostScript, Sun Raster, ASCII
Slide Shows	PostScript, GIF, HTML, Java
Tools	Solaris object code

Table 1: Data formats supported by Taxila.

File format conversion is a very important consideration for maintaining the versatility of the digital library.

WWW browsers such as Netscape Navigator which display HTML documents will normally allow GIF/TIFF files as inlined images. Other formats such as PostScript cause a new window to be displayed with the contents of the file displayed in the new window. Therefore, in order to store a research paper that contains embedded figures and tables, the following conversions are performed: (a) documents originally written in \LaTeX^2 are converted into HTML using the $\text{\LaTeX}2\text{HTML}$ conversion program, (b) figures included in such documents are converted from PostScript format into GIF format using a *ps2gif* convertor, and then made “transparent”, if required, using a *transgif* convertor which also translates the image from GIF87a format to GIF89a format, (c) tables are converted into HTML Table format (which can be displayed by Netscape Navigator version 3.0).

Among the representation schemes that we are investigating for future enhancements to the digital library is an object-oriented database model (using a database system such as Illustra) in which the information in the digital library will be stored as a hierarchical set of objects. Each object will contain the data associated with that object as well as functions that determine the category of the information and the tools required to present the information to the user. Thus, all the data conversion and interpretation operations for an object will be done by the functions associated with that object, thus making the interface transparent to not only the user but also to all the higher levels of objects in the library.

7. DOCUMENT BROWSER APPLICATION

Being a research center that conducts research in document analysis and recognition, we naturally incorporate document analysis and recognition capabilities into our digital library. We are conducting research into extracting relevant information from the scanned raster image of a document in response to user needs. The digital library contains many scanned papers that are stored in digitized raster form. The objective is to do some or all of the following on a given document image: (a) perform *layout analysis* on the document to find all the text and graphics regions; (b) label all the text regions according to their identities such as “title”, “footnote”, etc.; (c) identify basic logical linkages such as that between a photograph and its caption; (d) group related regions into *logical units* such as a “newspaper story”; (e) determine the *reading order* of the text blocks within a logical unit; and (f) selectively read the text in a given logical unit. We are developing a *document browser* capability into the digital library which will perform the above functions. This would enable the user to extract any relevant information while browsing a document.

A “document browser” application is being developed based on this approach. The browser allows a user to interactively specify queries on the documents in the digital library using a graphical user interface, provides feedback on the candidate documents and their properties at each stage of the query and retrieval process, and allows refinements of the query based on the intermediate results of the search. Results of a query can be displayed either as an image or as formatted text, and the browser allows the user to specify the level of detail in the displayed results.

8. RESULTS

Experimental results for this approach have been very encouraging. A large proportion of CEDAR researchers utilize the digital library to obtain information on techniques and related material relevant to their research. For example, researchers from various project groups have been using the digital library for information on handwritten word recognition, which is used in several of our current projects, as well as for archived material on OCR (optical character recognition) algorithms, which are used in CEDAR’s character recognition engines to read addresses on envelopes, names on forms, etc. As the material in the digital library has become more comprehensive, and its organization and retrieval capabilities have become more extensive, usage has grown to the point that the digital library has become the one centralized resource to which all CEDAR researchers can

turn for a variety of comprehensive, accurate and up-to-date information related to past and current CEDAR research.

The DeLoS system, which is used to identify parts of documents to be retrieved, has been trained using data from page images of various newspaper and journal pages, and extensively tested on data from digitized images of newspaper pages.

The DeLoS system has been tested on a variety of newspaper pages (e.g., *The Buffalo News*, *USA Today*). Overall, the DeLoS system performed fairly well for images from newspapers. Figure 2 shows the performance of the system for the *USA Today* newspaper pages, in terms of percentages of the original blocks correctly classified, grouped and read-ordered. It also shows the percentages of correctly segmented blocks that are correctly classified, grouped and read-ordered. Performance results for *The Buffalo News* followed a similar pattern (and have been described in detail in⁵). As we can see from the table in Figure 2, block segmentation and block type-categorization in the original images proved to be a deciding factor in the performance of the system.

Document ID	Block Classif.	Block Grouping	Read-Ordering
Page01	96.1 %	73 %	72.7 %
Page02	96.7 %	80.6 %	92.3 %
Page03	85.1 %	100 %	100 %
Page04	84.2 %	73.6 %	87.5 %
Page05	91.3 %	82.6 %	71.4 %
Page06	85.1 %	66.6 %	80 %
Page07	88.4 %	84.6 %	100 %
Page08	100 %	94.4 %	100 %
Page09	70.3 %	85.1 %	57.1 %
Page11	100 %	97.1 %	86.6 %
Page12	82.6 %	69.5 %	83.3 %
Page13	84 %	48 %	75 %

BLOCK CLASSIF. = % of blocks correctly classified
BLOCK GROUPING = % of blocks correctly grouped
READ-ORDERING = % of text blocks correctly read-ordered

Figure 2: Performance of DeLoS on pages of *USA Today*.

9. CONCLUSIONS

We have presented an approach to retrieving desired components of a document in response to user queries. The structure of a document structure analysis system (DeLoS) which enables this, was also presented. This approach allows a user to retrieve logical components of a document (e.g., a newspaper article, a citation list, a photograph, a chart, etc.). It also provides the ability to perform multi-stage index building based on the query, so that only those documents in the library that are candidates for satisfying a given query have to be completely analyzed, thus saving considerable search time.

We have also presented the design and implementation of Taxila, the CEDAR digital library for research on document analysis and recognition. This digital library is being used extensively by CEDAR researchers to

retrieve material of interest.

The CEDAR digital library is a prototype library for multimedia research information, and can be used as a model for WWW-based information repositories. It is possible for heterogeneous research information sources in non-local Web sites to be linked to the Taxila structure, thereby contributing to more efficient and effective sharing of research information among researchers at distributed locations. In addition, Taxila acts as an interactive document analysis and recognition tool for researchers, so that parts of a document can be selectively analyzed, and this not only furthers the state-of-the-art in this dynamic research area, but also provides a logical linkage between the areas of document understanding and digital library technology. Thus, the Taxila should provide a resource for those interested in the processes of creating digital document libraries as well as automating data entry from paper.

10. REFERENCES

1. S.W. Lam and S.N. Srihari. Multi-domain document layout understanding. In *Proceedings of ICDAR-91*, pages 112–120, Saint-Malo, France, Sept. 30–Oct. 2, 1991.
2. L. Lamport. *LaTeX: A Document Preparation System*. Addison-Wesley, 1986.
3. M. Landy, Y. Cohen, and G. Sperling. HIPS: A Unix-based image processing system. *Computer Vision, Graphics, and Image Processing*, 25:331–347, 1984.
4. A.M. Nazif and M.D. Levine. Low level image segmentation: An expert system. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6(5):555–577, September 1984.
5. D. Niyogi. *A Knowledge-Based Approach to Deriving Logical Structure from Document Images*. PhD thesis, State University of New York at Buffalo, 1994.
6. D. Niyogi and S.N. Srihari. An integrated approach to document decomposition and structural analysis. *International Journal of Imaging Systems and Technology*, 7:330–342, 1996.
7. D. Niyogi and S.N. Srihari. A rule-based system for document understanding. In *Proceedings of AAAI-86*, volume 2, pages 789–793, Philadelphia, PA, August 15–22, 1986.
8. D. Niyogi and S.N. Srihari. Knowledge-based derivation of document logical structure. In *Proceedings of ICDAR '95 (Third International Conference on Document Analysis and Recognition)*, Montreal, Canada, August 1995.
9. D. Niyogi, S.N. Srihari, and V. Govindaraju. Analysis of printed forms. In H. Bunke and P.S.P. Wang, editors, *Handbook on Optical Character Recognition and Document Image Analysis*. World Scientific Publishing Co., Singapore, 1996.
10. R.K. Srihari. PICTION: A system that uses captions to label human faces in newspaper photographs. In *Proceedings of AAAI-91*, pages 80–85, Anaheim, CA, 1991.
11. S.N. Srihari, S.W. Lam, J.J. Hull, R.K. Srihari, and V. Govindaraju. Intelligent data retrieval from raster images of documents. In *Proceedings of Digital Libraries '94 (The First Annual Conference on the Theory and Practice of Digital Libraries)*, College Station, TX, June 1994.