

# A Complex Document Information Processing Prototype

S. Argamon<sup>1</sup>, G. Agam<sup>1</sup>, O. Frieder<sup>1</sup>, D. Grossman<sup>1</sup>, D. Lewis<sup>2</sup>, G. Sohn<sup>3</sup>, K. Voorhees<sup>3</sup>

<sup>1</sup>Dept. of Computer Science, Illinois Institute of Technology, Chicago, IL, USA

<sup>2</sup>David D. Lewis Consulting, 858 W. Armitage Ave., #296, Chicago, IL, USA

<sup>3</sup>Clarabridge Inc., 3130 Fairview Park Drive, Falls Church, VA, USA

## ABSTRACT

We developed a prototype for integrated retrieval and aggregation of diverse information contained in scanned paper documents. Such *complex document information processing* combines several forms of image processing together with textual/linguistic processing to enable effective analysis of complex document collections, a necessity for a wide range of applications. This is the first system to attempt integrated retrieval from complex documents; we report its current capabilities.

## Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval

## Keywords

CDIP, Document image processing, Integrated retrieval

## 1. INTRODUCTION

*Complex Document Information Processing* (CDIP) involves the analysis of large masses of scanned paper documents, which often contain non-textual information such as handwriting, logos, or signatures. It is a critical problem in many key application areas, including litigation, intelligence analysis, knowledge management, and humanities scholarship. However, while specific solutions do exist for component problems such as OCR, handwriting analysis, logo recognition, and signature matching, no system to date has integrated extraction and analysis methods to enable users to pose queries that integrate these different forms of document information. Current practice, therefore, is to use separate processing systems independently, so that collating and cross-checking information items must be done by hand, underscoring the need for integrated systems. Furthermore, results can be degraded when individual processing modules are unaware of the larger context in which they operate; an integrated system will likely do better.

We describe a first research prototype for integrated CDIP, to enable work on these important research problems (see also [4]). Elsewhere [2] we describe our work on building a testbed collection

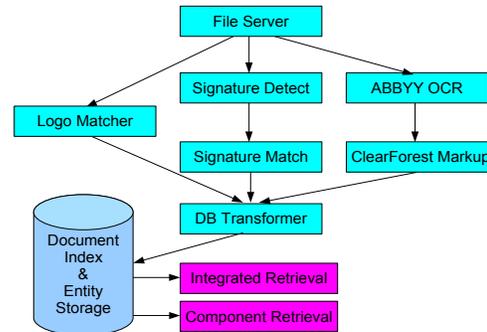


Figure 1: Architecture of the current prototype.

for CDIP, based on document images from the Legacy Tobacco Documents Library (<http://legacy.library.ucsf.edu/>).

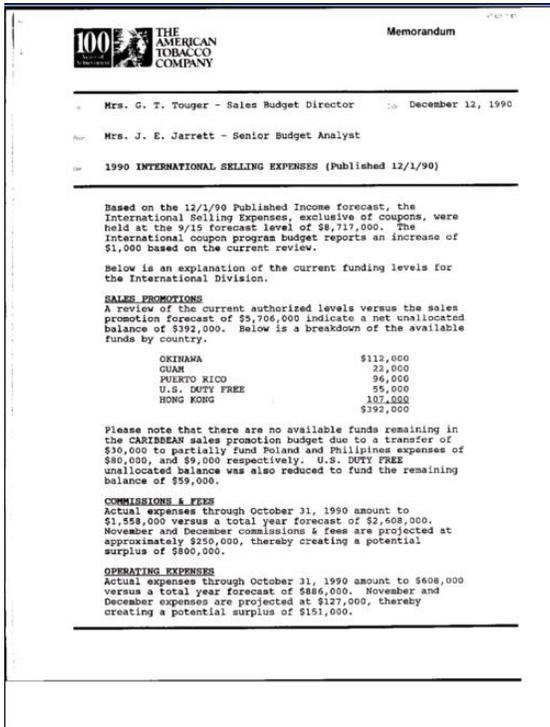
The prototype currently contains modules that extract and analyze text, signatures, and logos from complex documents, enabling integrated document image retrieval. The system (Fig. 1) contains a set of ingestion modules, which process different forms of document image data, importing them into a database schema which includes traditional document indices by keyword and relations between documents and their components (logos, signatures, and named entities). Retrieval operates via SQL queries on this unified database.

Each component in the figure is a separate thread, so that processing is fully parallelized and pipelined. Image files are served to processing modules dealing different types of document image information. (Future development will add preprocessing for noise removal, skew-correction, orientation determination, and region zoning.) The ABBYY OCR engine ([www.abbyy.com](http://www.abbyy.com)) is used to extract text from the document image. This text is fed to the ClearForest ([www.clearforest.com](http://www.clearforest.com)) information extraction module, which finds and classifies various named entities and relations. Signatures are segmented<sup>1</sup> and then fed to CEDAR's signature recognition system [3] which matches document signatures to known signatures in a database. Logos are segmented and matched using the DocLib package [1]. These three threaded processing paths are then synchronized, and the data extracted are transformed into a unified database schema for retrieval and analysis.

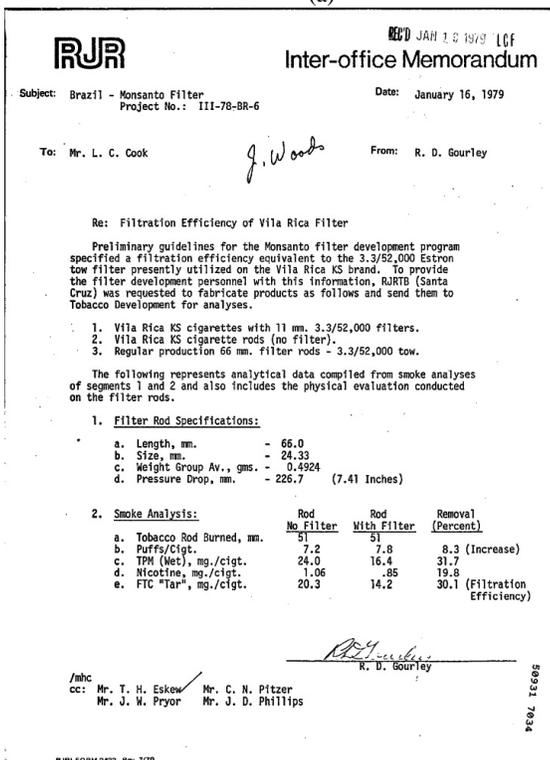
## 2. RETRIEVAL EXAMPLES

The rich collection of attributes our system associates with each document (including words, linguistic entities such as names and

<sup>1</sup>Signature segmentation is not yet integrated, and is currently simulated by hand-segmentation and labeling.



(a)



(b)

Figure 2: Retrieval results. (a) With ATC logo, the words “income forecast”, and mention of than \$500,000. (b) With R.J.Reynolds logo, the words “filtration efficiency”, and a signature.

Signature Group	Me	Ent Dollars
	Gertenbach, RF	\$37,454,447.88
	Schan, M	\$30,885,327.00
	Boffa, JR	\$17,420,705.00
	Nielsen, VG	\$958,354.82
	Bergman, JI	\$635,397.44

Figure 3: Prototype results showing signatures associated with the most total dollars (see text).

amounts, logos, and signatures) enables both novel forms of text retrieval, and the evidence combining capabilities of a relational database. While we have, as yet, no quantitative evaluations to report, we give examples here of the kinds of capabilities that our prototype currently supports. The mini-corpus used for this consists of 800 documents taken from the testbed we are building.

We consider integrated queries that our prototype makes possible for the first time. We apply conjunctive constraints on document image components to a straightforward document ranking based on total query-word frequency in the OCR'd document text; in Figure 2 we show document images retrieved for two such queries. The first is the unique document found containing both of the words “income forecast” as well as the American Tobacco Company logo and a dollar amount (a recognized entity type) greater than \$500K. The second example is the top-ranked document for “filtration efficiency” that also has the R.J. Reynolds logo and a signature. Note that neither of these documents would have been found just based on their printed text, as neither contains the company name explicitly.

In Figure 3 we show a ranked retrieval results for a document component query which asks for the five signatures with the highest total of dollar amounts mentioned in documents with each signature. This shows another novel way of integrating useful information extracted from document images which is easily implementable in our framework.

*Acknowledgment.* Supported in part by an ARDA Challenge Grant.

### 3. REFERENCES

- [1] K. Chen, S. Jaeger, G. Zhu, and D. Doermann. DOCLIB: A document processing research tool. In *SDIUT*, 2005.
- [2] D. D. Lewis, S. Argamon, G. Agam, O. Frieder, D. Grossman, and J. Heard. Building a test collection for complex document information processing. In *SIGIR-06*, 2006.
- [3] S. N. Srihari, M. K. Kalera, and A. Xu. Offline signature verification and identification using distance statistics. *IJPRAI*, 18(7), 2004.
- [4] S. S. Stein, S. Argamon, and O. Frieder. The effect of OCR on stylistic text classification. In *SIGIR-06*, 2006.