

# Robust Document Image Understanding Technologies

Henry Baird, Daniel Lopresti, Brian Davison, William Pottenger

Department of Computer Science and Engineering

Lehigh University

19 Memorial Drive West

Bethlehem, PA 18015

{baird|lopresti|davison|billp}@cse.lehigh.edu

## Abstract

*No existing document image understanding technology, whether experimental or commercially available, can guarantee high accuracy across the full range of documents of interest to industrial and government agency users. Ideally, users should be able to search, access, examine, and navigate among document images as effectively as they can among encoded data files, using familiar interfaces and tools as fully as possible. We are investigating novel algorithms and software tools at the frontiers of document image analysis, information retrieval, text mining, and visualization that will assist in the full integration of such documents into collections of textual document images as well as “born digital” documents. Our approaches emphasize versatility first: that is, methods which work reliably across the broadest possible range of documents.*

## 1 Introduction

The challenges faced by many industries and US government agencies in automating the capture, understanding, and reuse of scanned hardcopy documents include extremely high volumes of documents and dauntingly wide variety of document types. High-accuracy OCR systems do not exist for many languages and writing systems due to the lack of commercial incentives to develop them. Also, many documents, when scanned, yield images of such low quality that conventional OCR systems fail almost completely. Moreover, later-stage processes, including retrieval and data mining, may be severely impacted by document analysis and OCR errors.

In the Department of Computer Science and Engineering at Lehigh, we are studying many of these key issues. For example, in the past we have performed work on:

- high-accuracy OCR on low-quality document images [8, 44],
- robust retrieval from noisy text corpora by combining approximate string matching techniques with fuzzy logic [30, 28],
- document image quality modeling and applications of such models to the construction of high-performance OCR systems [1, 19],
- duplicate detection for scanned documents that have been subjected either to OCR [31] or character shape coding [29],
- the impact of recognition errors on document summarization [21], and
- automatic creation of hypertext links in document images, especially using figure references and bibliographic citations.

Recently, we identified some of the most pressing issues confronting government agencies attempting to build and manage large collections of scanned document images [2].

## 2 Research Directions

In this section, we discuss some of the research topics that we believe would help solve these problems and that we are capable of addressing.

### 2.1 “Versatility-First” DIA Research

One promising strategy for improving the performance of image understanding systems by the orders of magnitude that are needed is, we believe, to aim for *versatility first*. For decades the machine vision R&D community has optimized for high speed, and for high accuracy on some (often only a small) fraction of the input images, but only later — if at all — for versatility, by which we mean *guaranteed competence over a broad and precisely specified class of images*. As a result, vision technologies still fall

---

<sup>0</sup>Published in *Proceedings, 1st ACM Hardcopy Document Processing Workshop (HDP 2004)*, Washington, DC, November 12, 2004.

far short of both human abilities and users' needs: they are overspecialized, brittle, unreliable, and improving only with painful slowness.

A versatility-first vision research program begins when we select a broad, challenging family of images: *e.g.* all printed documents potentially containing any of many languages, scripts, page layout styles, and image qualities. Then, we investigate ways to:

- capture as much as possible of these images' variety in a formal generative (often stochastic) model that combines several submodels, *e.g.* of image quality, layout, and language (this requires both analytical rigor and sophisticated statistical modeling, for significant progress towards this, cf. [23, 2]);
- develop methods for inferring the parameters of such models from labeled training data (can be difficult even though there is a large relevant literature, *e.g.* [24, 37]);
- design provably optimal recognition algorithms, for each submodel, and for the system as a whole, for best possible results w.r.t. the models (an intellectual challenge but sometimes doable, *e.g.* [38, 39]);
- (only then) reduce runtimes to practical levels, carefully without loss of generality (this may require inventions but is almost always possible, *e.g.* [34, 7, 8]);
- organize the system to adapt its model parameters to unlabeled test data, on the fly, and so retrain itself with a minimum of manual assistance (progress has been reported, in recent years, at RPI [43], Bell Labs [5], and PARC [9]); and
- construct 'anytime' recognition systems which, when allowed to run indefinitely, are guaranteed to improve accuracy monotonically to the best achievable, *i.e.* consistent with the Bayes error of the problem (a daunting, exciting, as yet almost untouched research domain).

Our experience inventing, building, testing, patenting, and applying systems of this type has convinced us of their promise — successes so far include:

- a world record in accuracy (99.995% characters correct) achieved by exploiting semantic as well as syntactic models of image content (w/ Ken Thompson) [18];
- a page reader that is quickly and easily 'retargetable' to new languages including Japanese, Bulgarian, and Tibetan (w/ David Ittner, Tin Ho, & others) [3];
- an automatically self-correcting classifier that cuts its own error rate by large factors without retraining, given merely a single hint (w/ George Nagy) [5];
- a high-accuracy tabular-data reader that, with only 15 minute's clerical effort, can be trained to a new table-type, applied to over 400 different forms (w/ Tom Wood & John Shamilian) [22];

- a printed-text recognition technology, trainable with low manual effort, that maintains uniformly high accuracy over an unprecedentedly broad range of image qualities (w/ Gary Kopec & Prateek Sarkar, see refs above); and
- world-class web security technology (CAPTCHAs) able to block programs ('bots, spiders, etc) from abusing web services, by means of automated Turing tests that exploit the gap in ability between humans and machines in reading degraded images of text (w/ Allison Coates, Richard Fateman, Monica Chew *et al*) [4, 12].

## 2.2 Retrieving from Noisy Sources

Most published methods for retrieval of document images first attempt recognition and transcription followed by indexing and search operating on the resulting (in general, erroneous) encoded text (using, *e.g.*, standard "bag-of-words" information retrieval (IR) methods). Early papers by Taghva, *et al.* show that moderate error rates have little impact on the effectiveness of traditional information retrieval measures for relatively long documents [45, 46]. The excellent survey by Doermann [15] summarized the state of the art (in 1997) of retrieval of entire multi-page articles as follows:

1. at OCR character error rates below 5%, these IR methods suffer little loss of either recall or precision; and
2. at error rates above 20%, both recall and precision degrade significantly.

A crucial open problem, which we are studying, is the effectiveness of "first OCR, then IR" methods on short passages such as, in an extreme but practically important case, fields containing key metadata (such as title, author, etc). Approximate string matching techniques offer some promise for improving recall and precision [30, 28], as does the small but interesting literature on word-spotting in the image domain [40]. Within short passages of metadata, especially for old works, dictionary solutions may not help interpretation of arcane/rare/unique words (such as names of people, places).

## 2.3 Summarizing Noisy Documents

In a recent paper [21], we examined some of the challenges in summarizing noisy documents. In particular, we broke down the summarization process into four steps: sentence boundary detection, pre-processing (part-of-speech tagging [32] and syntactic parsing), extraction, and editing [20]. We tested each step on noisy documents and analyzed the errors that arose, finding that these modules suffered significant degradation as the noise level in the document increased. We also studied how the overall

quality of summarization was affected by the noise level and the errors made at each stage of processing.

In examining the accuracy of the OCR process using edit distance techniques [17], we determined that OCR performance varied widely depending on the type of degradation. Punctuation symbols were particularly hard-hit due to their small size, which is critical because of their importance in delimiting sentence boundaries. For clean text, sentence boundary detection is not a big problem; the reported accuracy is usually above 95% [36, 41, 42]. However, since such systems typically depend on punctuation, capitalization, and words immediately preceding and following punctuation to make judgments about potential sentence boundaries, detecting sentence boundaries in noisy documents is a challenge due to the unreliability of such features.

We also found that syntactic parsers may be very vulnerable to noise in a document. Even low levels of noise tended to lead to a significant drop in performance. For documents with high levels of noise, it may be better not to rely on syntactic parsing at all since it will likely fail on a large portion of the text, and even when results are returned, they will be unreliable.

Employing three measures used in the Document Understanding Conference [16] for assessing the quality of generated summaries, unigram overlap between the automatic summary and the human-created summary, bigram overlap, and the simple cosine, we evaluated the overall performance of our test summarization system. Not surprisingly, summaries of noisier documents generally had a lower overlap with human-created summaries (for full details, see [21]).

As our results showed, the methods we tested at every step were fragile, susceptible to failures and errors even with slight increases in the noise level of a document. Clearly, much work needs to be done to achieve acceptable performance in noisy document summarization. We need to develop summarization algorithms that do not suffer significant degradation when used on noisy documents. We also need to develop the robust natural language processing techniques that are required by summarization. These would include, for example, sentence boundary detection systems that can reliably identify sentence breaks in noisy documents.

## 2.4 Preserving Uncertainty

We would, in fact, like to preserve uncertainty throughout our system as much as possible. Doing so allows us to recognize where the system knows about possible errors, permitting better debugging, and possible incorporation of end-user correction and training. Given appropriate feedback about new

content, recognition systems can be trained, thus improving their performance on similar future tasks. In general, manual correction of OCR'd text is infeasible for large-scale efforts – the time and expense are too high. Instead, we propose to design and build a collaborative tool for editing and correction, providing valuable feedback to the underlying recognition model, both to train the system for future recognition tasks, but also to re-evaluate past uncertainty. Thus, the correction of one image from one page of a document could have a ripple effect throughout the document, and perhaps to other documents which had similar uncertainties.

Such a system will require work in a number of areas.

- Collaborative editing and community approval scheme, *a la* slashdot. A good editor will make corrections that are approved by others, increasing the editor's authority, thus decreasing the amount of confirmation required by others in the future.
- A strong dependency error model, so that when corrections are made, other scenarios with the same uncertainty can be quickly identified. Perhaps we'll need to go further – incorporating dependency information into all recognized text, not just uncertain text.

## 2.5 Automating Meta-Data Creation

A large portion of the expense and effort in bringing document images online is the creation of meta-data, including correcting OCR errors that may have arisen. In the case of Lehigh University's "Digital Bridges" digital library [14], the librarians involved in the project estimate that complete correction took approximately 10 minutes per page, or six pages an hour; so for a 300 page book, a total of 50 hours was required [33]. Based on feedback we have received, there is no doubt that the need for extensive manual post-processing is regarded as a major hurdle in the construction of large collections from scanned document images.

## 2.6 Automated Creation of Hypertext Links for Document Images

Text that references or discusses figures or other documents is an excellent candidate for innovative linking and navigation. We plan to identify, extract, and index such text, in addition to recognizing and indexing *within-image* text and explicit captions. This allows us to make non-text images (e.g., figures, plates) retrievable using text queries (in contrast to most content-based retrieval techniques [47, 49]).

Recognizing textual content that discusses a figure or image would also be useful in deciding to include an image for automated summarization purposes, and finding the first such reference can assist in re-flowing a document for better presentation (cf. [10]).

Prior work has focused on the automatic recognition and extraction of scholarly citations (*e.g.*, CiteSeer/ResearchIndex [25, 27]) but has not incorporated the discussion text as part of the cited document. On the Web, in contrast, search engines routinely associate the content in links both to the source document and to the cited Web page [11], since such text is a good descriptor of the target document [13]. This is exploited by motivated Web authors for search engine manipulation and for what is known as Google-bombing [6] — creating enough links with common anchor text to a particular site to place that site at or near the top of the rankings when the anchor text is used as a query.

When indexing images, the major Web search engines use some available text. They all use text within the URL of the image, but some go further. Google’s image search is capable of using image captions; it also uses page form text (pull-down menus), but not general text (or titles, etc.). AltaVista’s image search, in contrast, uses text from the citing page, which allows for many more matches, but also includes many poor matches.

The accurate selection of relevant text will make an otherwise unretrievable figure accessible. We plan to make use of text mining techniques [35]; in particular, given labeled examples of the kinds of references we wish to find, information extraction algorithms can be trained to recognize new occurrences (*e.g.*, as in [48]).

## 2.7 Recognizing Citations

Prior work has focused on recognition of citations in general (*e.g.*, CiteSeer/ResearchIndex) [27, 26] but has not incorporated the discussion text as part of the cited document. In contrast, on the Web, search engines routinely associate the content in links both to the source document and to the cited Web page, leading to current issues of Google-bombing [6]. On the Web, all resources are independent; instead, we are promoting an image from within a document to become as accessible as the text within a document. Google’s image search is capable of using image captions; also uses page form text (pull-down menus), but not general text (or titles, etc.) Recognizing textual content that discusses a figure or image would also be useful in deciding to include an image for automated summarization purposes (such as a news story, ala Google News). Recognizing the first reference to a figure is also important when re-flowing a

document for better presentation.

## 3 Conclusions

We have touched on several key areas of our research agenda. In the final camera-ready copy, we will report on progress towards these goals.

## References

- [1] H. Baird. Model-directed document image analysis. In *Proceedings of the DOD-sponsored Symposium on Document Image Understanding Technology (SDIUT 1999)*, pages 42–49, Annapolis, Maryland, April 1999.
- [2] Henry Baird, Kris Popat, Thomas Breuel, Prateek Sarkar, and Daniel Lopresti. Assuring high-accuracy document understanding: Retargeting, scaling up, and adapting. In *Proceedings of the Symposium on Document Image Understanding Technology*, pages 17–29, Greenbelt, MD, April 2003.
- [3] Henry S. Baird. Anatomy of a versatile page reader. *Proceedings of the IEEE*, 80(7):1059–1065, July 1992.
- [4] Henry S. Baird, Allison L. Coates, and Richard Fateman. Pessimism: a reverse turing test. *Int’l J. on Document Analysis and Recognition*, 5:158–163, 2003.
- [5] Henry S. Baird and George Nagy. A self-correcting 100-font classifier. In *Proceedings, IS&T/SPIE Symposium on Electronic Imaging: Science & Technology*, pages 106–115, San Jose, CA, February 1994.
- [6] ‘Miserable failure’ links to Bush: George W Bush has been Google bombed. BBC News, December 2003. <http://news.bbc.co.uk/2/hi/americas/3298443.stm>.
- [7] D. Bloomberg, T. Minka, and K. Popat. Document image decoding using iterated complete path search with subsampled heuristic scoring. In *Proceedings of the IAPR 2001 International Conference Document Analysis and Recognition (ICDAR 2001)*, Seattle, WA, September 2001.
- [8] T. Breuel and K. Popat. Recent work in the document image decoding group at xerox parc. In *Proceedings of the DOD-sponsored Symposium on Document Image Understanding Technology (SDIUT 2001)*, Columbia, Maryland, April 2001.
- [9] T. M. Breuel. Modeling the Sample Distribution for Clustering OCR. In *SPIE Conference on Document Recognition and Retrieval VIII*, 2001.
- [10] Thomas M. Breuel, William C. Janssen, Kris Popat, and Henry S. Baird. Paper to pda. In *Proc., IAPR 16th ICPR*, pages Vol. 4, 476–479, Quebec City, Canada, August 2002.

- [11] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual Web search engine. In *Proceedings of the Seventh International World Wide Web Conference*, Brisbane, Australia, April 1998.
- [12] M. Chew and H. S. Baird. Baffletext: a human interactive proof. In *Proc., 10th IS&T/SPIE Document Recognition & Retrieval Conf.*, Santa Clara, CA, January 2003.
- [13] Brian D. Davison. Topical locality in the Web. In *Proceedings of the 23rd Annual ACM International Conference on Research and Development in Information Retrieval (SIGIR 2000)*, pages 272–279, Athens, Greece, July 2000.
- [14] Digital Bridges. Lehigh University Libraries, April 2004. <http://bridges.lib.lehigh.edu/>.
- [15] D. Doermann. The indexing and retrieval of document images: A survey. *Computer Vision and Image Understanding*, 70(3), June 1998. Special Issue on “Document Image Understanding and Retrieval,” J. Kanai and H. S. Baird (Eds.).
- [16] Document Understanding Conference (DUC): Workshop on Text Summarization, 2002. <http://tides.nist.gov/>.
- [17] Jeffrey Esakov, Daniel P. Lopresti, and Jonathan S. Sandberg. Classification and distribution of optical character recognition errors. In *Proceedings of Document Recognition I (IS&T/SPIE Electronic Imaging)*, volume 2181, pages 204–216, San Jose, CA, February 1994.
- [18] Ken Thompson Henry S. Baird. Reading chess. *PAMI*, 12(6):552–559, 1990.
- [19] Tin Kam Ho and Henry S. Baird. Large-scale simulation studies in image pattern recognition. *IEEE Trans. on PAMI*, 19(10):1067–1079, October 1997.
- [20] Hongyan Jing. *Cut-and-paste Text Summarization*. PhD thesis, Department of Computer Science, Columbia University, New York, NY, 2001.
- [21] Hongyan Jing, Daniel Lopresti, and Chilin Shih. Summarizing noisy documents. In *Proceedings of the Symposium on Document Image Understanding Technology*, pages 111–119, Greenbelt, MD, April 2003.
- [22] Thomas L. Wood John H. Shamilian, Henry S. Baird. A retargetable table reader. pages 1–14, Ulm, Germany, August 1997.
- [23] T Kanungo, H Baird, and R Haralick. Estimation and validation of document degradation models. In *Proc. 4th Annual Symposium on Document Analysis and Information Retrieval*, Las Vegas, NV, April 1995.
- [24] G. Kopec. An em algorithm for character template estimation. submitted March 1997; returned for revision, but not revised due to the author’s death; available from PARC by request.
- [25] Steve Lawrence, Kurt Bollacker, and C. Lee Giles. Autonomous citation matching. In Oren Etzioni, editor, *Proceedings of the Third International Conference on Autonomous Agents*, New York, 1999. ACM Press.
- [26] Steve Lawrence, Kurt Bollacker, and C. Lee Giles. Distributed error correction. In *Digital Libraries 99 - The Fourth ACM Conference on Digital Libraries*, New York, 1999. ACM Press.
- [27] Steve Lawrence, C. Lee Giles, and Kurt Bollacker. Digital libraries and Autonomous Citation Indexing. *IEEE Computer*, 32(6):67–71, 1999.
- [28] Daniel Lopresti. Robust retrieval of noisy text. In *Proceedings of the Third Forum on Research and Advances in Digital Libraries*, pages 76–85, Washington, DC, May 1996.
- [29] Daniel Lopresti and A. Lawrence Spitz. Comparing the utility of optical character recognition and character shape coding in duplicate document detection. In *Proceedings of the Fourth IAPR International Workshop on Document Analysis Systems*, pages 439–450, Rio de Janeiro, Brazil, December 2000.
- [30] Daniel Lopresti and Jiangying Zhou. Retrieval strategies for noisy text. In *Proceedings of the Fifth Annual Symposium on Document Analysis and Information Retrieval*, pages 255–269, Las Vegas, NV, April 1996.
- [31] Daniel P. Lopresti. A comparison of text-based methods for detecting duplication in scanned document databases. *Information Retrieval*, 4(2):153–173, July 2001.
- [32] M. McCord. *English Slot Grammar*. IBM, 1990.
- [33] Philip Metzger. Private communication, April 2004.
- [34] Thomas P. Minka, Dan S. Bloomberg, and Kris Popat. Document image decoding using iterated complete path heuristic. In *Proceedings of IS&T/SPIE Electronic Imaging 2001: Document Recognition and Retrieval VIII*, San Jose, CA, January 2001.
- [35] S. Y. Mironova, M. W. Berry, S. Atchley, M. Beck, T. Wu, L. E. Holzman, W. M. Pottinger, and D. J. Phelps. Advancements in text mining algorithms and software, 2003.
- [36] D. Palmer and M. Hearst. Adaptive multilingual sentence boundary disambiguation. *Computational Linguistics*, 23(2):241–267, June 1997.
- [37] Kris Popat. Decoding of text lines in grayscale document images. In *Proceedings of the 2001*

- International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2001)*, Salt Lake City, Utah, May 2001. IEEE. To appear.
- [38] Kris Popat, Dan Bloomberg, and Dan Greene. Adding linguistic constraints to document image decoding. In *Proc., 4th International Workshop on Document Analysis Systems*, Rio de Janeiro, Brazil, December 2000. International Association of Pattern Recognition.
- [39] Kris Popat, Dan Greene, Justin Romberg, and Dan S. Bloomberg. Adding linguistic constraints to document image decoding: Comparing the iterated complete path and stack algorithms. In *Proceedings of IS&T/SPIE Electronic Imaging 2001: Document Recognition and Retrieval VIII*, San Jose, CA, January 2001.
- [40] Toni M. Rath and R. Manmatha. Features for word spotting in historical manuscripts. In *Proceedings of the Seventh International Conference on Document Analysis and Recognition*, pages 218–222, August 2003.
- [41] J. C. Reyner and A. Ratnaparkhi. A maximum entropy approach to identifying sentence boundaries. In *Proceedings of the Fifth Conference on Applied Natural Language Processing*, Washington D.C., 1997.
- [42] M. Riley. Some applications of tree-based modelling to speech and language. In *Proceedings of the DARPA Speech and Natural Language Workshop*, pages 339–352, Cape Cod, MA, 1989.
- [43] P Sarkar. *Style Consistency in Pattern Fields*. PhD thesis, Rensselaer Polytechnic Institute, May 2000.
- [44] P. Sarkar, H. S. Baird, and X. Zhang. Training on severely degraded text–line images. [submitted to] IAPR Int’l Conf. on Document Analysis & Recognition, Edinburgh, August, 2003.
- [45] Kazem Taghva, Julie Borsack, and Allen Condit. Effects of OCR errors on ranking and feedback using the vector space model. *Information Processing and Management*, 32(3):317–327, 1996.
- [46] Kazem Taghva, Julie Borsack, and Allen Condit. Evaluation of model-based retrieval effectiveness with OCR text. *ACM Transactions on Information Systems*, 14:64–93, January 1996.
- [47] R.C. Veltkamp and M. Tanase. Content-based image retrieval systems: A survey. Technical Report UU-CS-2000-34, Department of Computing Science, Utrecht University, October 2000. <http://citeseer.ist.psu.edu/veltkamp00contentbased.html>.
- [48] Tianhao Wu and William M. Pottenger. A semi-supervised active learning algorithm for information extraction from textual data. *Journal of the American Society for Information Science and Technology*, 2004. In press.
- [49] A. Yoshitaka and T. Ichikawa. A survey on content-based retrieval for multimedia databases. *IEEE Transactions on Knowledge and Data Engineering*, 11(1), Jan.–Feb. 1999.