

REPRESENTING OCRED DOCUMENTS IN HTML

Tao Hong and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo
Buffalo, New York 14228
email: {taohong,srihari}@cedar.buffalo.edu

ABSTRACT: OCR is an error-prone process. It is time-consuming and expensive to manually proofread OCR results. The errors remaining in OCRed texts can cause serious problems in reading and understanding if they do not refer to the original image representation. As demonstrated in this paper, a hybrid document which combines symbolic representation and image representation may relieve the problem. If we represent a OCRed document properly in HTML, OCR errors will not have much negative effect on the human reading process in a HTML browser and can be corrected by using a HTML authoring tool. Under the approach, an experiment evaluating a Japanese OCR system developed in CEDAR is also reported in this paper.

1 Overview of the Approach

OCR is a process to transform a given document from its image representation into its symbolic representation. After this process, we obtain a text document which is electronically searchable, indexable and reusable. However, the transformation is error-prone, especially when the image quality is poor. Therefore, if we just save the pure text document generated by the OCR process, the text may have some potential errors. Although errors in a OCRed text may not cause problem in tasks such as text categorization and retrieval, they may have a negative effect if the text is to be read by a human being.

In a current OCR system, although its final output usually is just the symbolic content, the detailed information on segmentation and recognition is stored as an intermediate step. The detailed information usually includes: (1) coordinates of the bounding boxes of text blocks, lines, words and characters; (2) alternatives and confidence scores of recognition for each character/word image. Based on this information, many potential OCR errors or suspicious points can be marked. Traditionally, a OCR system provides a specially designed graphical user interface (GUI) for an end-user to proofread OCRed text and manually correct those errors. Figure 1 shows a OCRed text generated by a commercial OCR system.

It is very time-consuming and expensive to manually proofread OCR results if a large number of scanned documents have to be transformed into texts. Therefore, for OCRed documents in large digital libraries, usually no proofreading is applied. By reading a OCRed text as shown in Figure 1, people may feel frustrated because OCR errors prevent them from getting exact information from the document.

Although it was originally designed for presenting documents on WWW, HTML is rapidly developing into the most popular document representation and a universal, platform-independent client interface for a broad range of applications running on Internet/Intranet and network-centric application servers. Because HTML allows embedded images inside normal text, a OCRed document can be represented in HTML as a combination of text and images. The point is to make potential OCR errors transparent to our reading process: for those words which are recognized with high confidence, recognition results are used; otherwise, original word images are placed. Figure 2 shows the HTML file for the OCRed text in Figure 1. Similarly, Figure 3 illustrates a small Japanese document image and its OCR result in HTML. Although there are more than 6,000 Kanji characters in the Japanese character set, the Japanese OCR system used here is designed to recognize about 3,000 JIS level 1 Kanji characters which are often used in Japan. When characters from JIS level 2 are encountered, recognition results for those characters will always be incorrect. It is best to keep those characters as original images so that the text will still be

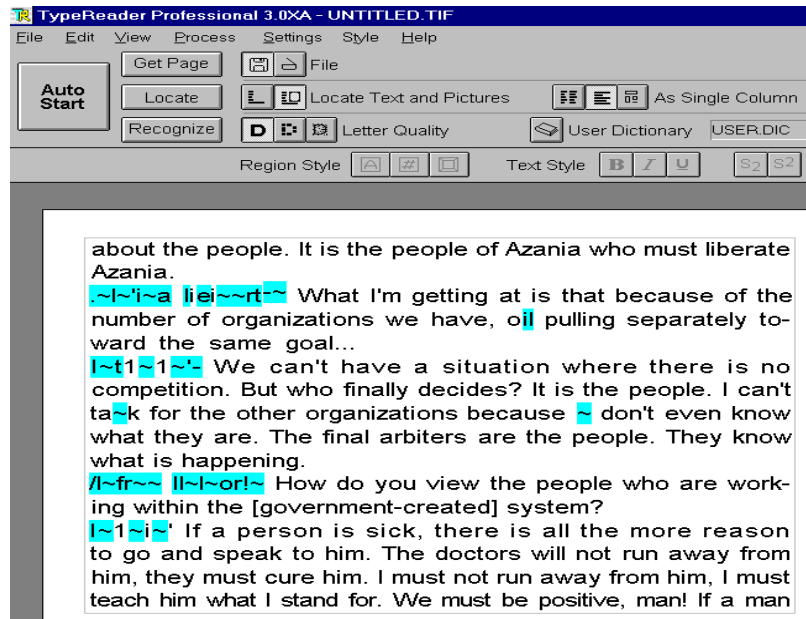


Figure 1: The OCR result of a text block generated by a commercial OCR system, TypeReader 3.0 from ExperVision Inc.. In the graphical user interface for proofreading, the places where the system failed or had low confidence are highlighted. An user can correct those errors after examining the original image.

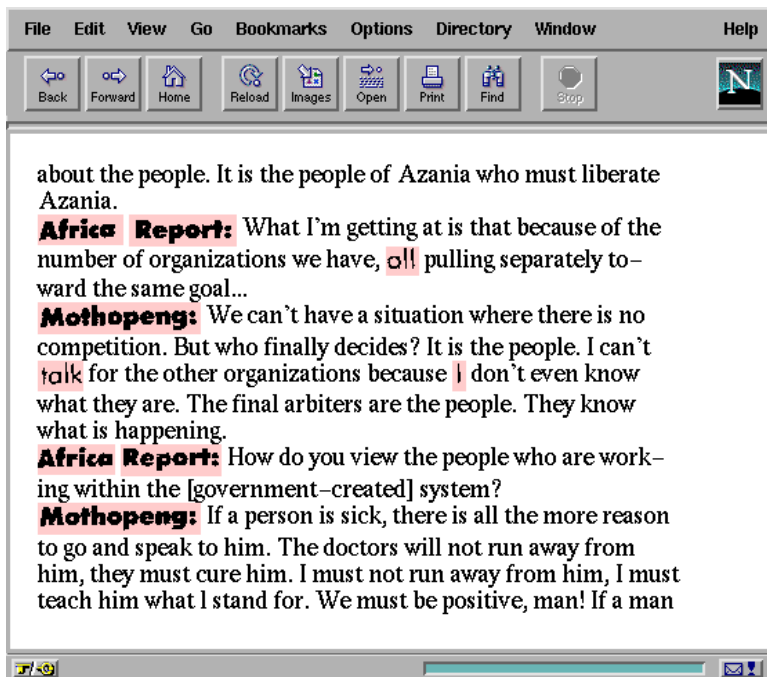
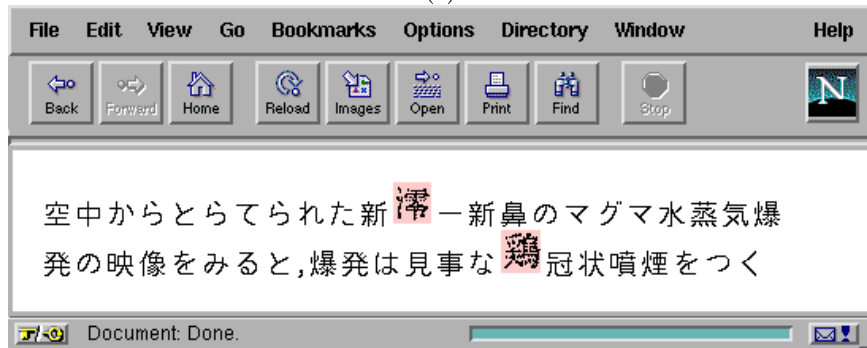


Figure 2: The HTML file generated from the OCRed text in Figure 1 is displayed in Netscape Navigator. It is a hybrid representation of text and images. Word images from the original document image are embedded when the OCR system failed or had low confidence in recognizing the word.

空中からとらえられた新澤一新鼻のマグマ水蒸気爆
発の映像をみると、爆発は見事な鶏冠状噴煙をつく

(a)



(b)

Figure 3: An example of a Japanese document image (a) and its OCR result in HTML (b). The Japanese OCR engine can only recognize about 3,000 often used Kanji characters (the first level of JIS code). Two Kanji characters above have no way to be recognized correctly because they are from the second level of JIS code. They are displayed as images in the HTML file because of recognition failure.

readable.

Because HTML allows hyperlinks between different entities, in a HTML file for a OCRred document, we can create hyperlinks between a word string and its image so that a user can check the original image whenever he/she feels uncertain about the recognition result. We will demonstrate this later in our experiment.

Hybrid representation of OCRred documents is not a new concept. In Adobe's Acrobat Capture [1], it has already been used to encode a recognized text into a PDF file which is much compressed and can be rendered the same as the original image in Acrobat Reader. Here, we propose to use HTML for the hybrid representation. The main advantages of doing so are: HTML is a simple but powerful markup language for the network environment, and there are many HTML browsing and authoring tools available on different platforms.

In the process of developing an experimental Japanese OCR system, we realized that a Web browser such as Netscape Navigator can be a nice environment for OCR system evaluation. By automatically converting recognition results into hyperlinked HTML files, we can do OCR result evaluation, verification and editing easily on different platforms across a local network or on the Internet. By taking this approach, much programming work on GUI design and multilingual support can be avoided.

2 An Experiment on Evaluating Japanese Document Recognition Using HTML

Cherry Blossom is a prototype system for Japanese document recognition which has been developed at CEDAR in the past few years [2]. This system is designed to meet several challenges of machine-printed Japanese document recognition: wide variations in data quality (facsimile, photocopy, newspaper etc.),

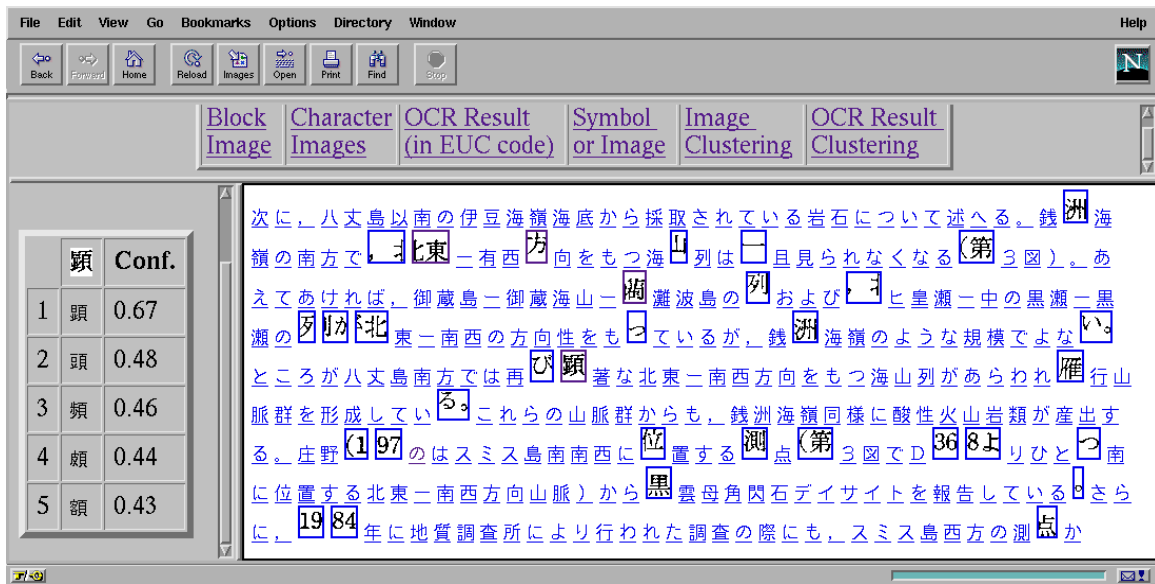


Figure 4: Evaluating OCR results using HTML and Netscape Navigator. There are three frames in the homepage. The top one is the TOC frame which has a table in which each entry is a hyperlink to one of six HTML files to be displayed in the right bottom frame, the MAIN frame. The left bottom one is the OCR frame which is used to display the recognition result for a character image. The HTML file currently being displayed in the MAIN frame is the hybrid representation of the recognition results, in which each character position is either a JIS code or an image. Each character position in the Main frame is also a hyperlink to its OCR result. By clicking on a character position, its OCR result will be displayed in the OCR frame.

low scan resolution (e.g., 200ppi), large character set (about 3,000 Kanji characters in the first level JIS code), and a variety of font and point sizes. Given a Japanese document page, the system can deskew the image, segment the page into blocks, discriminate between text/non-text blocks, determine the alignment of a text block, segment text regions into lines and further into character images, and recognize character images.

In order for rapid prototyping of evaluation tools which are flexible and platform-dependent, a framework for integrating the JOCR system with Internet/Intranet has been outlined. Under the framework, the output of the JOCR system has to be converted to HTML files, which can be viewed by using a Web browser such as Netscape Navigator. Simple Perl scripts can be written to convert the OCR results into HTML files. More detail description of the experiment can be found in [3].

For a Japanese document page, the recognition system can generate OCR results for each text block. In the current experiment, we focused on the task of converting OCR results for a text block into HTML files. The aspects of system performance to be evaluated are: character segmentation, character recognition and character clustering¹. Given OCR results of a text block, a Perl script was written to generate six HTML files. In each of those HTML files, for every character position (either a JIS code or an image), a hyperlink to its OCR file is provided so that its OCR result can be examined by clicking the mouse on it. The OCR file of a character is also in HTML. In the file, the character image is extracted from the block image and saved in GIF format; the classification result is also saved in a HTML file.

¹The HTML demonstrations for various Japanese and Chinese document pages are located at <http://www.cedar.buffalo.edu/JOCR/HTML-demo/index.html>.

The Netscape Navigator browser provides a mechanism to have several frames simultaneously inside its window. Each frame can display an individual HTML file. Using the browser, a homepage for JOCR results is designed with three frames: the TOC frame, the OCR frame and the MAIN frame.

Figure 4 is a snapshot of a homepage automatically generated for a OCRred Japanese document. Activated by selecting the entry `Symbol_or_Image` in the TOC frame, a hybrid HTML file for OCR results can be displayed in the MAIN frame. Here, each character is either a JIS code or an image. For this text block, most characters are in JIS because their OCR confidence scores are high. For those characters with low OCR confidence, their original images are displayed. Obviously, those positions where images are placed usually have segmentation problems or recognition errors. Although there are many segmentation and recognition errors in the OCR result which makes the result in JIS code difficult to read, the content in the hybrid representation is quite readable and not much information has been lost due to segmentation and recognition errors. For several positions, OCR results are incorrect but with high confidence. Therefore, JIS codes are placed there. But it is not difficult for a user to detect those errors in their local context and by examining their original images. If the user wants to correct those errors, he can load the hybrid HTML file into a HTML editor and replace those images and wrong JIS codes with the correct JIS codes.

It is easy to modify the Perl script when the appearance of a HTML page needs to be changed. In our experiment, the performance of character segmentation, character recognition and image-based clustering were evaluated. New scripts for converting OCR results to HTML files can be designed rapidly when some other aspects of OCR performance need to be studied.

3 Summary

OCR is an error-prone process. The errors remaining in the OCRred text can cause serious problems in reading and understanding. As demonstrated in this paper, a hybrid document which combines symbolic representation and image representation may relieve the problem. Such a hybrid document in HTML is searchable and also error-recoverable. If we represent a OCRred document in such a hybrid format, OCR errors will not have much negative effect on the human reading process. Therefore, the hybrid HTML can be recommended as an ideal format for OCRred documents in an online digital library [4].

References

- [1] <http://www.adobe.com/Acrobat/Capture.html>. Adobe's Acrobat Capture.
- [2] S. N. Srihari, G. Srikantan, T. Hong, and B. Grom. A general-purpose Japanese optical character recognition system. In *Symposium on Document Analysis and Information Retrieval*, pages 271–286, April 15-17 1996.
- [3] T. Hong, S.F. Wu, and S.N. Srihari. Evaluating japanese document recognition in the internet/intranet environment. In *IAPR Workshop On Document Analysis Systems*, pages 633–650, Oct. 14-16 1996.
- [4] M. Lesk. Substituting images for books: The economics for libraries. In *Symposium on Document Analysis and Information Retrieval*, pages 1–16, April 15-17 1996.