

# A Unified Approach Towards Text Recognition

Tao Hong and Jonathan J. Hull\* and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition  
Department of Computer Science  
State University of New York at Buffalo  
Buffalo, New York 14260  
{taohong,srihari}@cedar.buffalo.edu

\*RICOH California Research Center  
2882 Sand Hill Road, Suite 115  
Menlo Park, CA 94025  
hull@crc.ricoh.com

## ABSTRACT

In our recent research, we found that visual inter-word relations can be useful for different stages of English text recognition such as character segmentation and postprocessing. Different methods had been designed for different stages. In this paper, we propose a unified approach to use visual contextual information for text recognition. Each word image has a lattice, which is a data structure to keep results of segmentation, recognition and visual inter-word relation analysis. A lattice allows ambiguity and uncertainty at different levels. A lattice-based unification algorithm is proposed to analyze information in the lattices of two or more visually related word images, and upgrade their contents. Under the approach, different stages of text recognition can be accomplished by the same set of operations—inter-word relation analysis and lattice-based unification. The segmentation and recognition result of a word image can be propagated to those visually related word images and can contribute to the recognition of them. In this paper, the formal definition of lattice, the unification operators and their uses are discussed in detail.

Key Words: OCR, text recognition, visual inter-word constraints, character segmentation, postprocessing, and unification.

## 1 Introduction

Character segmentation, recognition and postprocessing are three major components of an OCR system. Traditionally, they are treated as very different tasks and organized as a cascade (see Figure 1.(a)). For example, character segmentation is an image processing task, which tries to segment each word image into a sequence of character images; character recognition is a pattern

classification task, which aims at mapping each character image onto its symbolic identification; and postprocessing is a symbolic processing task, which utilizes contextual information, such as lexical knowledge, to detect and correct potential character recognition errors.

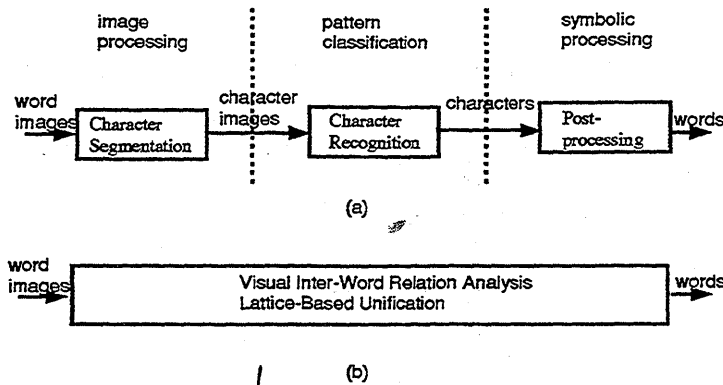


Figure 1: Models of text recognition. (a) the traditional approach; (b) the approach pursued in this paper.

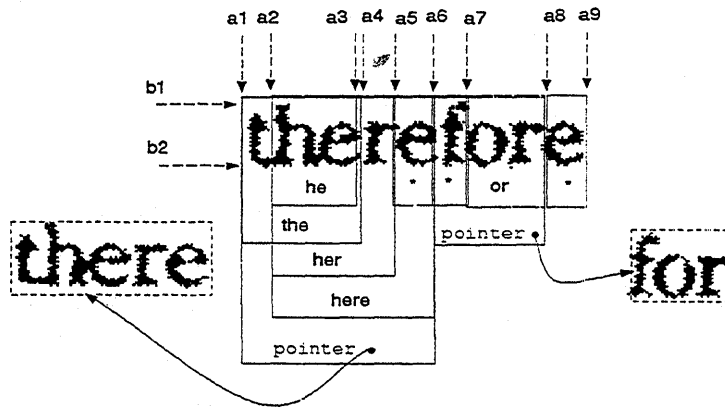
For a degraded text, character segmentation and recognition usually are error-prone. Lattice has been used to keep alternatives of character segmentation points and character recognition results [1, 2, 4, 3, 12, 13, 14]. In previous studies, each entry in a lattice is usually limited to a simple representation at the character level.

Our recent work demonstrated that visual context can be useful for text recognition [8, 9, 10, 11]. We designed different algorithms to use visual inter-word relations for character segmentation and postprocessing. In this paper, we extend our previous work and propose a unified approach for text recognition by exploiting visual inter-word context. Under the approach, different stages of text recognition can be accomplished by the same set of operations—inter-word relation analysis and lattice-based unification (see Figure 1.(b)).

Given a text image, we assume that word segmentation has been done [5]. Each word image can be represented by the coordinates of its bounding box in the text image. For each word image, a lattice can be defined as a data structure which keeps record of the current result for the word image. The result can be derived from character segmentation, character recognition, or visual context analysis. An entry in a lattice is a sub-image associated with its label which reflects recognition result for the piece of image. A label can be a character, a string or a pointer to another image. At beginning, the lattice for each word image is empty. New entries can be created and added to a lattice by using traditional methods of character segmentation and recognition, or, as we propose here, through a unification algorithm. The detail of the approach, including the unification algorithm, will be described in following sections.

## 2 Definition

Given a text image, an image  $I$  corresponding to a word or a string can be represented by its bounding box,  $(s_x, s_y, e_x, e_y)$ , where  $(s_x, s_y)$  is the coordinate of the left-top pixel and  $(e_x, e_y)$  the one of the right-bottom pixel of the image. Here, we assume that the word image has been de-slanted if it was originally printed in an italic font. A coordinate value  $x$  here can approximately represent any coordinate in the range  $[x - \delta, x + \delta]$  where  $\delta$  is a small integer.



$$\begin{aligned}
 &L(a1, b1, a9, b2) \\
 &= \{(a1, b1, a4, b2, \text{"the"}, 0.99), (a1, b1, a6, b2, \text{pointer}, 0.90), \\
 &\quad (a2, b1, a3, b2, \text{"he"}, 0.90), (a2, b1, a5, b2, \text{"her"}, 0.90), \\
 &\quad (a2, b1, a6, b2, \text{"here"}, 0.95), (a5, b1, a6, b2, \text{"*"}, 0.90), \\
 &\quad (a6, b1, a7, b2, \text{"*"}, 0.95), (a6, b1, a8, b2, \text{pointer}, 0.90), \\
 &\quad (a7, b1, a8, b2, \text{"or"}, 0.88), (a8, b1, a9, b2, \text{"*"}, 0.90)\}
 \end{aligned}$$

Figure 2: The lattice for a word image

For an image  $I(s_x, s_y, e_x, e_y)$ , its lattice can be defined as a set.

$$\begin{aligned}
 &L(s_x, s_y, e_x, e_y) \\
 &= \{(u, s_y, v, e_y, \text{label}, \text{conf}) \mid s_x \leq u < v \leq e_x\} \quad (1)
 \end{aligned}$$

where  $(u, s_y, v, e_y, \text{label}, \text{conf})$  is the entry for a fragment of the image  $I(s_x, s_y, e_x, e_y)$ :  $(u, s_y, v, e_y)$  is its bounding-box;  $\text{label}$  is its label; and  $\text{conf}$  is the confidence for the label. A label can be defined as

$$\text{label} = * \mid \text{character} \mid \text{string} \mid \text{pointer}$$

where "\*" is a special label which indicates no recognition result yet for the segment of the image, and *pointer* is a pointer to the bounding box of another image fragment which is visually equivalent to the segment. For the same image segment, there can be several entries with different labels.

A lattice can store results of segmentation, recognition or visual similarity analysis. For example, Figure 2 shows a word image, therefore, and its lattice. There are ten entries in the lattice. Those entries with the label \* are segmentation results; those ones with string labels are recognition results; and those ones with the label pointer are results of visual inter-word relation analysis. A pointer entry links to another word/sub-word image which is visually equivalent to the image fragment represented by the entry.

According to the definition, a lattice for a large image can be decomposed as the union of the lattices corresponding to the fragments of the image,

$$\begin{aligned}
 L(s_x, s_y, e_x, e_y) &= \{(u, s_y, v, e_y, label, conf) \mid s_x \leq u < v \leq e_x\} \\
 &= \bigcup_{s_x \leq u < v \leq e_x} L(u, s_y, v, e_y)
 \end{aligned} \tag{2}$$

A word lattice,  $L(s_x, s_y, e_x, e_y)$ , can be viewed as a directed graph, in which, each entry  $(u, s_y, v, e_y, label, conf)$  is an edge from the vertex  $u$  to the vertex  $v$ ; the start vertex is  $s_x$ ; and the end vertex is  $e_x$ .

Given a word lattice,  $L(s_x, s_y, e_x, e_y)$ , its interpretation can be defined as a set of complete paths from  $s_x$  to  $e_x$ ,

$$\{(s_x, s_y, e_x, e_y, label\_of\_the\_path, conf\_of\_the\_path) \mid label\_of\_the\_path \in \text{DICTIONARY or } conf\_of\_the\_path \text{ is high}\} \tag{3}$$

Complete paths can be computed by using a dynamic programming algorithm [6]. If two paths  $(a1, s_y, a2, e_y, label_1, conf_1)$  and  $(a2, s_y, a3, e_y, label_2, conf_2)$  have been found, a longer path

$$(a1, s_y, a3, e_y, label_1 \bullet label_2, \frac{(a2 - a1) * conf_1 + (a3 - a2) * conf_2}{a3 - a1})$$

can be defined. Here, " $\bullet$ " denotes "concatenation".

Figure 3 (b) shows an example how two complete paths can be derived through the search on the lattice in Figure 3 (a). Two word candidates in Figure 3 (c) are generated as word recognition result for the word image.

We also define  $I(s1_x, s1_y, e1_x, e1_y) \approx I(s2_x, s2_y, e2_x, e2_y)$  if the two images can match each other very well. Details of visual similarity measurement and inter-word relation calculation can be found in our previous papers [8, 9, 10, 11].

### 3 Lattice-Based Unification

Given two images  $I_1 = I(s1_x, s1_y, e1_x, e1_y)$  and  $I_2 = I(s2_x, s2_y, e2_x, e2_y)$ , if  $I_1 \approx I_2$ , the lattice  $L(s1_x, s1_y, e1_x, e1_y)$  for  $I_1$  can be upgraded to a new lattice,  $L(s1_x, s1_y, e1_x, e1_y) \uplus$

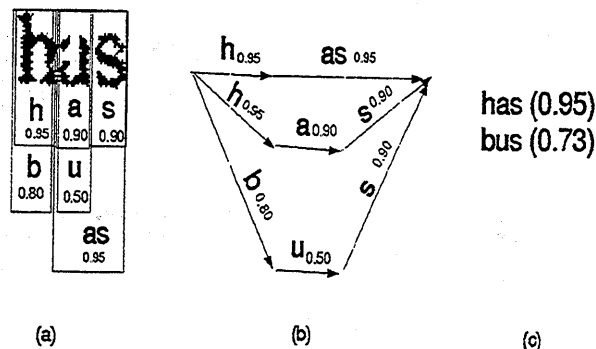


Figure 3: Generating word recognition result from a lattice. (a) the lattice; (b) complete paths; and (c) word candidates

$L(s2_x, s2_y, e2_x, e2_y)$ , which is the unification of the lattice  $L(s1_x, s1_y, e1_x, e1_y)$  and the lattice  $L(s2_x, s2_y, e2_x, e2_y)$ . Formally,

$$\begin{aligned}
 & L(s1_x, s1_y, e1_x, e1_y) \uplus L(s2_x, s2_y, e2_x, e2_y) \\
 &= L(s1_x, s1_y, e1_x, e1_y) \cup \\
 & \quad \{(u, s1_y, v, e1_y, label, conf) \mid s1_x \leq u < v \leq e1_x \text{ and} \\
 & \quad (s2_x + u - s1_x, s2_y, s2_x + v - s1_x, e2_y, label, conf) \in L(s2_x, s2_y, e2_x, e2_y)\} \quad (4)
 \end{aligned}$$

Similarly,  $L(s2_x, s2_y, e2_x, e2_y)$  can be updated to

$$\begin{aligned}
 & L(s2_x, s2_y, e2_x, e2_y) \uplus L(s1_x, s1_y, e1_x, e1_y) \\
 &= L(s2_x, s2_y, e2_x, e2_y) \cup \\
 & \quad \{(u, s2_y, v, e2_y, label, conf) \mid s2_x \leq u < v \leq e2_x \text{ and} \\
 & \quad (s1_x + u - s2_x, s1_y, s1_x + v - s2_x, e1_y, label, conf) \in L(s1_x, s1_y, e1_x, e1_y)\} \quad (5)
 \end{aligned}$$

Figure 4 shows the effect of the unification operation on two visually equivalent word images. For each individual image, its lattice keeps results of character segmentation and recognition. As you can see, there are some segmentation and recognition errors for each individual image. After unification of those two lattices, the lattice for each image is expanded and the correct path - Africa, can be generated from it.

If the image  $I_1 = I(s1_x, s1_y, e1_x, e1_y)$  matches a part of the image  $I_2 = I(s2_x, s2_y, e2_x, e2_y)$ ,  $I(s1_x, s1_y, e1_x, e1_y) \approx I(m1_x, s2_y, m2_x, e2_y)$  where  $s2_x \leq m1_x < m2_x \leq e2_x$ , we can define

$$\begin{aligned}
 & L(s1_x, s1_y, e1_x, e1_y) \uplus L(s2_x, s2_y, e2_x, e2_y) \\
 &= L(s1_x, s1_y, e1_x, e1_y) \cup \\
 & \quad \{(u, s1_y, v, e1_y, label, conf) \mid s1_x \leq u < v \leq e1_x \text{ and} \\
 & \quad (m1_x + u - s1_x, s2_y, m1_x + v - s1_x, e2_y, label, conf) \in L(s2_x, s2_y, e2_x, e2_y)\} \quad (6)
 \end{aligned}$$

Similarly,

$$L(s2_x, s2_y, e2_x, e2_y) \uplus L(s1_x, s1_y, e1_x, e1_y)$$

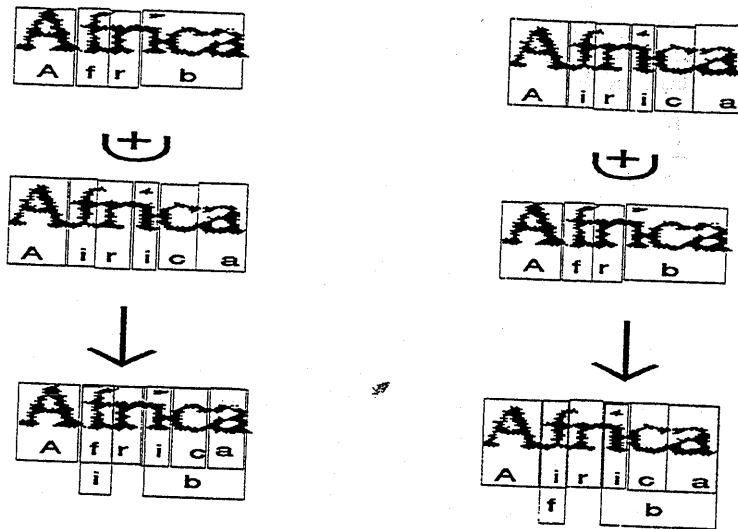


Figure 4: The unification of two visually equivalent images

$$\begin{aligned}
 &= L(s2_x, s2_y, e2_x, e2_y) \cup \\
 &\quad \{(m1_x, s2_y, m2_x, e2_y, "*", conf)\} \cup \\
 &\quad \{(u, s2_y, v, e2_y, label, conf) \mid s2_x \leq u < v \leq e2_x \text{ and} \\
 &\quad (s1_x + u - m1_x, s1_y, s1_x + v - m1_x, e1_y, label, conf) \in L(s1_x, s1_y, e1_x, e1_y)\} \quad (7)
 \end{aligned}$$

Examples of the unification operation defined above are illustrated in Figure 5.

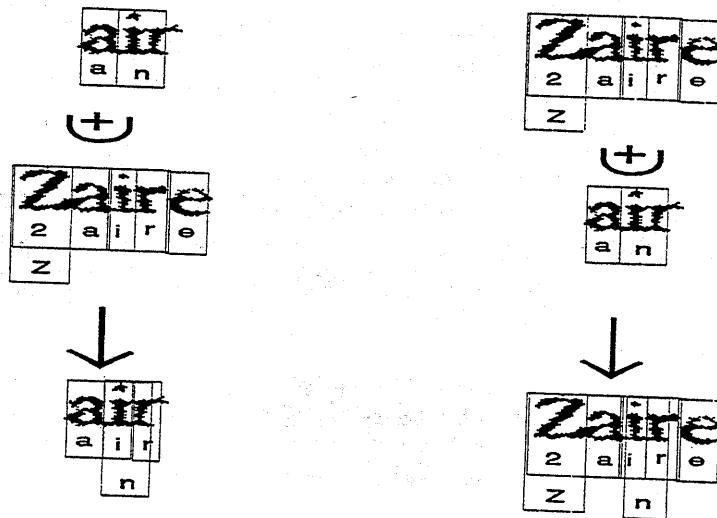


Figure 5: The unification of two word images

In general case, if two images,  $I(s1_x, s1_y, e1_x, e1_y)$  and  $I(s2_x, s2_y, e2_x, e2_y)$  can partially match,  $I(m1_x, s1_y, m2_x, e1_y) \approx I(n1_x, s2_y, n2_x, e2_y)$  where  $s2_x \leq m1_x < m2_x \leq e2_x$  and  $s2_x \leq$

$n1_x < n2_x \leq e2_x$ , the unification operations between the two lattices can be defined as

$$\begin{aligned}
 & L(s1_x, s1_y, e1_x, e1_y) \uplus L(s2_x, s2_y, e2_x, e2_y) \\
 &= L(s1_x, s1_y, e1_x, e1_y) \cup \\
 & \quad \{(u, s1_y, v, e1_y, label, conf) \mid m1_x \leq u < v \leq m2_x \text{ and} \\
 & \quad (n1_x + u - m1_x, s2_y, n1_x + v - m1_x, e2_y, label, conf) \in L(s2_x, s2_y, e2_x, e2_y)\} \quad (8)
 \end{aligned}$$

and

$$\begin{aligned}
 & L(s2_x, s2_y, e2_x, e2_y) \uplus L(s1_x, s1_y, e1_x, e1_y) \\
 &= L(s2_x, s2_y, e2_x, e2_y) \cup \\
 & \quad \{(u, s2_y, v, e2_y, label, conf) \mid n1_x \leq u < v \leq n2_x \text{ and} \\
 & \quad (m1_x + u - n1_x, s1_y, m1_x + v - n1_x, e1_y, label, conf) \in L(s1_x, s1_y, e1_x, e1_y)\} \quad (9)
 \end{aligned}$$

Figure 6 gives an example of the unification defined above.

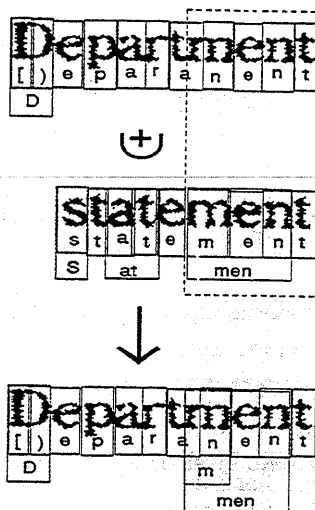


Figure 6: The unification of two word images

Given a lattice  $L(s_x, s_y, e_x, e_y)$  which is derived from a sequence of unification operations, some segmentation points may be represented implicitly in the lattice. To represent them explicitly, the lattice can be expanded by applying the following operation,

$$\begin{aligned}
 & \odot(L(s_x, s_y, e_x, e_y)) \\
 &= L(s_x, s_y, e_x, e_y) \cup \\
 & \quad \{(a1, s_y, a2, e_y, "*", conf_1 \times conf_2), \\
 & \quad (a2, s_y, a3, e_y, "*", conf_1 \times conf_2), \\
 & \quad (a3, s_y, a4, e_y, "*", conf_1 \times conf_2) \mid \\
 & \quad a1 \leq a2 < a3 \leq a4 \\
 & \quad \text{and } (a1, s_y, a3, e_y, *, conf_1) \in L
 \end{aligned}$$

$$\begin{aligned}
& \text{and } (a_2, s_y, a_4, e_y, *, \text{conf}_1) \in L \} \cup \\
& \{(a_1, s_y, a_2, e_y, "*", \text{conf}_1 \times \text{conf}_2), \\
& (a_3, s_y, a_4, e_y, "*", \text{conf}_1 \times \text{conf}_2) \mid \\
& a_1 \leq a_2 < a_3 \leq a_4 \\
& \text{and } (a_1, s_y, a_4, e_y, "*", \text{conf}_1) \in L \\
& \text{and } (a_2, s_y, a_3, e_y, "*", \text{conf}_1) \in L \} - \\
& \{(u, s_x, v, e_y, "*", \_)\} \mid \\
& (u, s_x, v, e_y, \text{label}, \_) \in L \text{ where } \text{label} \neq "*" \} \quad (10)
\end{aligned}$$

where “\_” means that anything can fill the field. Figure 7 shows the effect of this operation on an example image.

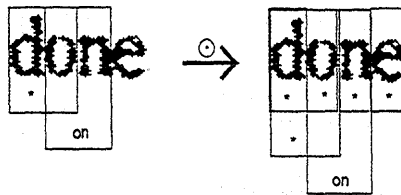


Figure 7: Re-organize a lattice

## 4 Text Recognition through Lattice-based Unification

Degraded text recognition still is a challenging task. In our previous studies [7, 8, 10, 11], we found:

- For English words, six types of inter-word relations can be defined to describe similarities among words, either at the symbolic level or at the image level;
- Relations at the symbolic level and at the image level are consistent in principle; visually equivalent images have to be interpreted/recognized consistently at the symbolic level;
- There are abundant inter-word relations in a normal English text. Visual inter-word relations can be calculated from text image by using simple image processing techniques, even when the image is highly degraded;
- Traditional methods of character segmentation, recognition and postprocessing often generate errors which violate consistency at the symbolic and image levels for visually equivalent word images;
- Character segmentation can be conducted by using visual inter-word relation; many OCR errors can be recovered by using visual inter-word relations.

The unification operations we defined above provide us a framework to do character segmentation, recognition and postprocessing uniformly, at least in principle. Figure 8 shows an



example that how the word image "dates" can be segmented into character level through two operations,  $\oplus$  and  $\ominus$ . In Figure 8 (a), an word image "at" is used to split the word image; in Figure 8 (b) the word image is further segmented by using another word image "date".

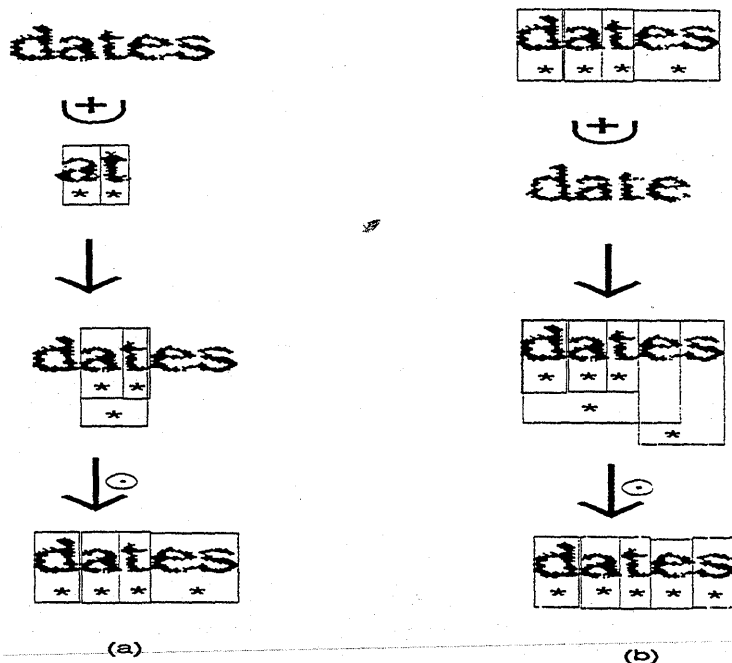


Figure 8: Character segmentation by unification

Give an entry  $(u, s_x, v, e_y, "*", -)$  in a lattice, an OCR can be applied to test if the image fragment can be recognized as a character with high confidence. If so, the character recognition result will replace the "\*" in the entry. Traditional segmentation algorithm can also be use to provide segmentation points for a lattice.

Although the operations in the paper are defined for two word images, it can be easily extended as operations on a image cluster which has more than two visually equivalent images.

## 5 Conclusions

In the paper we proposed a new approach for visual text recognition. Under the approach, visual inter-word relations are utilized uniformly to overcome difficulties of segmentation, recognition and postprocessing for individual word images. Unification operations were defined for the purpose and their uses for text recognition under the approach were discussed. Further experimental work is being conducted to show the approach is practical.

## References

- [1] H.S. Baird. Anatomy of a versatile page reader. *Proceedings of the IEEE*, 80(7):1059-1065, 1992.
- [2] R.G. Casey. Character segmentation in document ocr: Progress and hope. In *Symposium on Document Analysis and Information Retrievals*, pages 13-40, April 24-26 1995.
- [3] C.H. Chen. Lexicon-driven word recognition. In *Proceedings of the Third International Conference on Document Analysis (ICDAR-95)*, pages 919-922, August 14-16 1995.
- [4] C.H. Chen. Structuring a large lexicon for word recognition. In *Symposium on Document Analysis and Information Retrievals*, pages 359-366, April 24-26 1995.
- [5] S. Chen, R. M. Haralick, and I. T. Phillips. Simultaneous word segmentation from document images using recursive morphological closing transform. In *Proceedings of the Third International Conference on Document Analysis (ICDAR-95)*, pages 761-764, August 14-16 1995.
- [6] T.H. Cormen, C.E. Leiserson, and R.L. Rivest. *Introduction to Algorithms*. The MIT Press, 1990.
- [7] T. Hong. Integration of visual inter-word constraints and linguistic knowledge in degraded text recognition. In *Proceedings of 32nd Annual Meeting of Association for Computational Linguistics, Student Sessions*, pages 328-330, 27-30 June 1994.
- [8] T. Hong and J. J. Hull. Improving ocr performance with word image equivalence. In *Symposium on Document Analysis and Information Retrievals*, pages 177-190, April 24-26 1995.
- [9] T. Hong and J. J. Hull. Visual inter-word relations and their use for ocr postprocessing. In *Proceedings of the Third International Conference on Document Analysis (ICDAR-95)*, August 14-16 1995.
- [10] T. Hong and J. J. Hull. Visual inter-word relations and their use in character segmentation. In *Proceedings of the Conference on Document Recognition, 1995 SPIE Symposium (SPIE95)*, February 1995.
- [11] Tao Hong. *Degraded Text Recognition Using Visual and Linguistic Context*. PhD thesis, Computer Science Department of SUNY at Buffalo, 1995.
- [12] J. Schurmann, N. Bartneck, T. Bayer, J. Franke, E. Mandler, and M. Oberlander. Document analysis - from pixels to contents. *Proceedings of the IEEE*, 80(7):1101-1119, 1992.
- [13] C.C. Tappert, C.Y. Suen, and T. Wakahara. The state of the art in on-line handwriting recognition. *IEEE Transactions on pattern analysis and machine intelligence*, 12(8):787-808, 1990.
- [14] S. Tsujimoto and H. Asada. Resolving ambiguity in segmenting touching characters. In *1st Int. Conf. on Document Analysis and Recognition*, pages 701-709, October 1991.



SPIE—The International Society for Optical Engineering

PROCEEDINGS

# *Document Recognition III*

**Luc M. Vincent**  
**Jonathan J. Hull**  
*Chairs/Editors*

**29–30 January 1996**  
**San Jose, California**

*Sponsored by*  
IS&T—The Society for Imaging Science and Technology  
SPIE—The International Society for Optical Engineering

*Published by*  
SPIE—The International Society for Optical Engineering



**Volume 2660**

SPIE is an international technical society dedicated to advancing engineering and scientific applications of optical, photonic, imaging, electronic, and optoelectronic technologies.