Horst Bunke
A. Lawrence Spitz (Eds.)

# Document Analysis Systems VII

**7th International Workshop, DAS 2006**
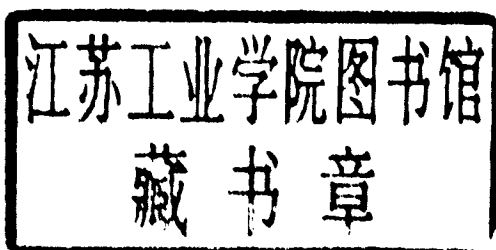**Nelson, New Zealand, February 2006**
**Proceedings**

IAPR

Springer

Horst Bunke   A. Lawrence Spitz (Eds.)

# Document Analysis Systems VII

7th International Workshop, DAS 2006
Nelson, New Zealand, February 13-15, 2006
Proceedings

Springer

Volume Editors

Horst Bunke
University of Bern
Department of Computer Science
Neubrückstr. 10, 3012 Bern, Switzerland
E-mail: bunke@iam.unibe.ch

A. Lawrence Spitz
DocRec Ltd
34 Strathaven Place, Atawhai, Nelson 7001, New Zealand
E-mail: spitz@docrec.com

# Lecture Notes in Computer Science 3872

# Lecture Notes in Computer Science

For information about Vols. 1–3771

please contact your bookseller or Springer

# Preface

DAS 2006 is the Seventh International Association for Pattern Recognition Workshop on Document Analysis Systems and was held in Nelson, New Zealand. DAS 2006 built on the tradition of past workshops held in Kaiserslautern, Germany (1994), Malvern, PA (1996), Nagano, Japan (1998), Rio de Janeiro, Brazil (2000), Princeton, NJ (2002), and Florence, Italy (2004). The goal of this meeting is to bring together those who have designed systems, or systems components, to solve real-world problems in document analysis.

Document analysis systems is inherently an interdisciplinary field encompassing such diverse disciplines as image processing, pattern recognition, document structure and natural language processing. DAS 2006 attempted to bring these disciplines together and to provide interactions between systems developers, suppliers and end users.

We received 78 papers from 19 countries. Each submission was reviewed by three reviewers. In addition to the Program Committee members, 42 other reviewers helped in this process. From those submissions and their reviews, we went through the difficult and sometimes painful process of ranking papers for acceptance or rejection. In the end we accepted 33 papers for oral presentation and 22 for presentation at poster sessions.

We, the Co-chairmen of DAS 2006, wish to express our gratitude to all of our colleagues who have reviewed the papers submitted for this conference.

We are proud to have brought two distinguished speakers to Nelson for keynote addresses: Ian Witten of the University of Waikato, the father of the New Zealand Digital Library, and James Fruchterman, a pioneer in modern commercial optical character recognition and currently CEO of Benetech.

We owe a special debt of gratitude to Marcus Liwicki of the University of Bern for his tireless work at maintaining the website, managing the flow of papers and reviews into the ConfMan system and assembling the proceedings for publication by Springer. He was ably assisted by Andreas Schlapbach.

We are fortunate that Siemens, Hitachi and Humanware provided DAS with financial support, and we thank them for doing so. Additionally, following the DAS tradition, the organizers of DAS 2004 have passed on the surplus from running that workshop for our use.

But ultimately it is the collection of authors who submitted papers to DAS to whom we owe the greatest gratitude. It is on them and their high-quality submissions that the success of DAS 2006 relies.

February 2006

Horst Bunke and Larry Spitz
Program Chairs
DAS 2006

# Organization

DAS 2006 was organized by DocRec Ltd.

## Executive Committee

Conference Chairs:        Larry Spitz (DocRec Ltd, New Zealand)
Horst Bunke (University of Bern, Switzerland)

## Program Committee

Apostolos Antonacopoulos (UK)
Henry Baird (USA)
Thomas Breuel (Germany)
Horst Bunke (Switzerland)
Andreas Dengel (Germany)
David Doermann (USA)
Andrew Downton (UK)
Michael Fairhurst (UK)
Hiromichi Fujisawa (Japan)
Venugopal Govindaraju (USA)
Tin Kam Ho (USA)
Jianying Hu (USA)
Rolf Ingold (Switzerland)
Rangachar Kasturi (USA)
Koichi Kise (Japan)
Seong-Whan Lee (Korea)
Daniel Lopresti (USA)
Raghavan Manmatha (USA)
Simone Marinai (Italy)
Udo Miletzki (Germany)
Yasuaki Nakano (Japan)
Larry Spitz (New Zealand)
Karl Tombre (France)

## Referees

Stefan Agne
Andrew Bagdanov
Ardhendu Behera
Koustav Bhattacharya
Alain Biem

Jean-Luc Blöchle
Matthew Boonstra
Jakob Brendel
Joshua Candamo
Farzin Deravi

Faisal Farooq
Gunnar Grimnes
Richard Guest
Sanaul Hoque
Gareth Howells
Jonathan Hull
Masakazu Iwamura
Stefan Jaeger
Thomas Kieninger
Malte Kiesel
Bertin Klein
Dar-Shyang Lee
Hansheng Lei
Jian Liang
Rainer Lindwurm
Vasant Manohar
Dalila Mekhaldi

David Mihalcik
Tristan Miller
Pranab Mohanty
Sunita Nayak
Shinichiro Omachi
Christoph Pesch
Maurizio Rigamonti
Thomas Roth-Berghofer
Sven Schwarz
Karthik Sridharan
Seiichi Uchida
Himanshu Vajaria
Ludger van Elst
Shankar Vembu
Alan Yang

## Sponsoring Institutions

Siemens AG, Munich, Germany
HumanWare Group, Christchurch, New Zealand
Hitachi Central Research Laboratory, Tokyo, Japan

## Scientific Sponsors

DocRec Ltd, Atawhai, Nelson, New Zealand
University of Bern, Switzerland
International Association for Pattern Recognition

# Table of Contents

## Session 1: Digital Libraries

## Session 2: Image Processing

## Session 3: Handwriting 1

## Session 4: Document Structure and Format

## Session 5: Tables

## Session 6: Handwriting 2

## Session 7: Language and Script Identification

## Session 9: Systems and Performance Evaluation

## Session 10: Retrieval and Segmentation

## Posters

# Retrieval from Document Image Collections

A. Balasubramanian, Million Meshesha, and C.V. Jawahar

Centre for Visual Information Technology,
International Institute of Information Technology,
Hyderabad - 500 032, India
jawahar@iiit.ac.in

**Abstract.** This paper presents a system for retrieval of relevant documents from large document image collections. We achieve effective search and retrieval from a large collection of printed document images by matching image features at word-level. For representations of the words, profile-based and shape-based features are employed. A novel DTW-based partial matching scheme is employed to take care of morphologically variant words. This is useful for grouping together similar words during the indexing process. The system supports cross-lingual search using OM-Trans transliteration and a dictionary-based approach. System-level issues for retrieval (eg. scalability, effective delivery etc.) are addressed in this paper.

## 1 Introduction

Large digital libraries, such as Digital Library of India (DLI) [1] are emerging for archiving large collection of printed and handwritten documents. The DLI aims at digitizing all literary, artistic, and scientific works of mankind so as to create better access to traditional materials, easier preservation, and make documents freely accessible to the global society. More than 25 scanning centers all over India are working on digitization of books and manuscripts. The mega scanning center we have, has around fifty scanners, each one of them capable lof scanning approximately 5000 pages in 8 hours. As on September 2005, close to 100 thousand books with 25 million pages were digitized and made available online by DLI (*http://dli.iiit.ac.in*) as document images.

Building an effective access to these document images requires designing a mechanism for effective search and retrieval of textual data from document image collections. Document image indexing and retrieval were studied with limited scope in literature [2]. Success of these procedures mainly depends on the performance of the OCRs, which convert the document images into text. Much of the data in DLI are in Indian languages. Searching in these document image collections based on content, is not presently possible. This is because OCRs are not yet able to successfully recognize printed texts in Indian languages. We need an alternate approach to access the content of these documents [3]. A promising alternate direction is to search for relevant documents in image space without any explicit recognition. We have been motivated by the successful attempts on

**Fig. 1.** Conceptual Diagram of the Searching Procedure from Multilingual Document Image Database. A Web Demo for the Above Procedure is Available Online at http://cvit.iiit.ac.in/wordsearch.html.

locating a specific word in handwritten English documents by matching image features for historical documents [4, 5].

We have already addressed algorithmic challenges for effective search in document images [6] . This paper describes the issues associated with the implementation of a scalable system for Indian language document images. A conceptual block diagram of our prototype system is shown in Figure 1. Our system accepts textual query from users. The textual query is first converted to an image by rendering, features are extracted from these images and then search is carried out for retrieval of relevant documents. Results of the search are pages from document image collections containing queried word sorted based on their relevance to the query.

## 2   Challenges in Design and Implementation of the System

Search and retrieval from large collection of document images is a challenging task, specially when there is no textual representation available. To design and implement a successful search engine in image domain, we need to address the following issues.

*Search in images:* Search in image space requires appropriate representational schemes and similarity measures. Success of content-based image retrieval(CBIR) schemes were limited by the diversity of the image collections. Digital libraries primarily archive text images, but of varying quality, script, style, size and font.

We need to come up with appropriate features and matching schemes, which can represent the content (text), while invariant to the popular variations.

*Degradations of documents:* Documents in digital libraries are extremely poor in quality. Popular artifacts in printed document images include (a) Excessive dusty noise, (b) Large ink-blobs joining disjoint characters or components, (c) Vertical cuts due to folding of the paper, (d) Cuts of arbitrary direction due to paper quality or foreign material, (e) Degradation of printed text due to the poor quality of paper and ink, (f) Floating ink from facing pages etc. We need to design an appropriate representational scheme and matching algorithm to accommodate the effect of degradation.

*Need for cross-lingual retrieval:* Document images in digital libraries are from diverse languages. Relevant documents that users need may be available in different languages. Most educated Indians can read more than one language. Hence, we need to design a mechanism that allows users to retrieve all documents related to their queries in any of the Indian languages.

*Computational speed:* Searching from large collection of document images pass through many steps: image processing, feature extraction, matching and retrieval of relevant documents. Each of these steps could be computationally expensive. In a typical book, there could be around 90,000 words and processing all of them online is practically impossible. We do all computationally expensive operations during the offline indexing (Section 4) and do minimal operations during online retrieval (Section 5).

*Indian languages:* Indian languages pose many additional challenges [7]. Some of these are: (i) lack of standard representation for the fonts and encoding, (ii) lack of support from operating system, browsers and keyboard, and (iii) lack of language processing routines. These issues add to the complexity of the design and implementation of a document image retrieval system.

## 3   Representation and Matching of Word Images

Word images extracted from documents in digital libraries are of varying quality, script, font, size and style. An effective representation of the word images will have to take care of these artifacts for successful searching and retrieval. We combined two categories of features to address these effects: word profiles and structural features. Explicit definitions of these features may be seen in [6].

*Word Profiles:* Profiles of the word provide a coarse way of representing a word image for matching. Profiles like upper word, lower word, projection and transition profiles are used here for word representation. Upper and lower word profiles capture part of the outlining shape of a word, while projection and transition profiles capture the distribution of ink along one of the two dimensions in a word image.