# Lecture Notes in Computer Science 3163

Simone Marinai   Andreas Dengel (Eds.)

# Document Analysis Systems VI

6th International Workshop, DAS 2004
Florence, Italy, September 8 - 10, 2004
Proceedings

Springer

Volume Editors

Simone Marinai
Università di Firenze, Dipartimento di Sistemi e Informatica
Via S. Marta, 3 - 50139 Firenze, Italy
E-mail: marinai@dsi.unifi.it

Andreas Dengel
German Research Center for Artificial Intelligence (DFKI)
P.O.Box 2080, 67608 Kaiserslautern, Germany
E-mail: Andreas.Dengel@dfki.de

# Preface

This volume contains papers selected for presentation at the 6th IAPR Workshop on Document Analysis Systems (DAS 2004) held during September 8–10, 2004 at the University of Florence, Italy. Several papers represent the state of the art in a broad range of "traditional" topics such as layout analysis, applications to graphics recognition, and handwritten documents. Other contributions address the description of complete working systems, which is one of the strengths of this workshop. Some papers extend the application domains to other media, like the processing of Internet documents.

The peculiarity of this 6th workshop was the large number of papers related to digital libraries and to the processing of historical documents, a taste which frequently requires the analysis of color documents. A total of 17 papers are associated with these topics, whereas two years ago (in DAS 2002) only a couple of papers dealt with these problems.

In our view there are three main reasons for this new wave in the DAS community. From the scientific point of view, several research fields reached a thorough knowledge of techniques and problems that can be effectively solved, and this expertise can now be applied to new domains. Another incentive has been provided by several research projects funded by the EC and the NSF on topics related to digital libraries. Last but not least, the organization of focused events, like the recent DIAL workshop chaired by Henry Baird and Venu Govindraju in Palo Alto (CA), had a strong impact on the definition of new research directions. However, it is indeed a lucky coincidence that this new trend in DAS research emerged in this edition organized in a town such as Florence, which keeps such an exceptional artistic and cultural heritage.

We received a total of 79 submissions from 19 countries, and we selected 31 oral presentations and 22 posters highlighted with short oral introductions. As a supplement to this proceedings, notes from the workshop discussions and other material related to presented papers will be posted on the DAS 2004 website: http://www.dsi.unifi.it/DAS04. Each paper was reviewed by three reviewers whom we would like to warmly thank here. We should mention the valuable support and hints provided by members of the Program Committee and past DAS chairs. We also wish to acknowledge the generosity of our sponsors: the International Association for Pattern Recognition, the University of Florence, the DFKI, ABBYY, Hitachi, and Siemens.

Special thanks are due to Alessio Ceroni, Cristina Dolfi, and Emanuele Marino for their invaluable contributions to the local organization.

June 2004                                                            Simone Marinai
                                                                     Andreas Dengel

# Organization

## Workshop Co-chairs

| | |
|---|---|
| Simone Marinai | University of Florence, Italy |
| Andreas Dengel | DFKI, Germany |

## Program Committee

| | |
|---|---|
| Apostolos Antonacopoulos | University of Liverpool, UK |
| Henry Baird | Lehigh University, USA |
| Francesca Cesarini | University of Florence, Italy |
| David Doermann | University of Maryland, USA |
| Andrew Downton | University of Essex, UK |
| Hiromichi Fujisawa | Hitachi Central Research Laboratory, Japan |
| Jianying Hu | IBM T.J. Watson Research Center, USA |
| Rolf Ingold | University of Fribourg, Switzerland |
| Ramanujan Kashi | Avaya Labs Research, USA |
| Koichi Kise | Osaka Prefecture University, Japan |
| Dan Lopresti | Lehigh University, USA |
| Donato Malerba | University of Bari, Italy |
| Udo Miletzki | Siemens Dematic, Germany |
| Yasuaki Nakano | Kyushu University, Japan |
| Lambert Schomaker | Rijksuniversiteit Groningen, The Netherlands |
| Giovanni Soda | University of Florence, Italy |
| Larry Spitz | Document Recognition Technologies, New Zealand |
| Karl Tombre | LORIA-INPL, France |
| Luc Vincent | LizardTech, USA |
| Marcel Worring | University of Amsterdam, The Netherlands |

## Additional Referees

| | | |
|---|---|---|
| Annalisa Appice | Dimosthenis Karatzas | T.R. Roth-Berghofer |
| Margherita Berardi | Michele Lapi | Jane Snowdon |
| Alain Biem | Larry O'Gorman | Salvatore Tabbone |
| Thomas Breuel | Huanfeng Ma | Yefeng Zheng |
| Michelangelo Ceci | Gérald Masini | Gary Zi |
| Philippe Dosch | Eugene Ratzlaff | |
| Stefan Jaeger | Maurizio Rigamonti | |

# Table of Contents

## Digital Libraries

## Historical Documents

## Layout Analysis

## Color Documents

## Handwritten Documents

## Graphics Recognition

# Internet Documents

# Document Analysis Systems

# Applications