# Lecture Notes in Computer Science 3411

*Commenced Publication in 1973*
Founding and Former Series Editors:
Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Sung Hyon Myaeng   Ming Zhou
Kam-Fai Wong   Hong-Jiang Zhang (Eds.)

# Information Retrieval Technology

Asia Information Retrieval Symposium, AIRS 2004
Beijing, China, October 18-20, 2004
Revised Selected Papers

Springer

Volume Editors

Sung Hyon Myaeng
Information and Communications University (ICU)
119 Munji-Ro, Yuseong-Gu, Daejeon, 305-714, South Korea
E-mail: myaeng@icu.ac.kr

Ming Zhou
Hong-Jiang Zhang
Microsoft Research Asia
5F, Beijing Sigma Center
No. 49 Zhichun Road Haidian District, Beijing 100080, China
E-mail: {mingzhou,hjzhang}@microsoft.com

Kam-Fai Wong
The Chinese University of Hong Kong
Shatin, N.T., Hong Kong, China
E-mail: kfwong@se.cuhk.edu.hk

# Preface

The Asia Information Retrieval Symposium (AIRS) was established by the Asian information retrieval community after the successful series of Information Retrieval with Asian Languages (IRAL) workshops held in six different locations in Asia, starting from 1996. While the IRAL workshops had their focus on information retrieval problems involving Asian languages, AIRS covers a wider scope of applications, systems, technologies and theory aspects of information retrieval in text, audio, image, video and multimedia data. This extension of the scope reflects and fosters increasing research activities in information retrieval in this region and the growing need for collaborations across subdisciplines.

We are very pleased to report that we saw a sharp increase in the number of submissions and their quality, compared to the IRAL workshops. We received 106 papers from nine countries in Asia and North America, from which 28 papers (26%) were presented in oral sessions and 38 papers in poster sessions (36%). It was a great challenge for the Program Committee to select the best among the excellent papers. The low acceptance rates witness the success of this year's conference.

After a long discussion between the AIRS 2004 Steering Committee and Springer, the publisher agreed to publish our proceedings in the Lecture Notes in Computer Science (LNCS) series, which is SCI-indexed. We feel that this strongly attests to the excellent quality of the papers.

The attendees were cordially invited to participate in and take advantage of all the technical programs at this conference. A tutorial was given on the first day to introduce the state of the art in Web mining, an important application of Web document retrieval. Two keynote speeches covered two main areas of the conference: video retrieval and language issues. There were a total of eight oral sessions run, with two in parallel at a time, and two poster/demo sessions.

The technical and social programs, which we are proud of, were made possible by the hard-working people behind the scenes. In addition to the Program Committee members, we are thankful to the Organizing Committee (Shao-Ping Ma and Jianfeng Gao, Co-chairs), Interactive Posters/Demo Chair (Gary G. Lee), and the Special Session and Tutorials Chair (Wei-Ying Ma). We also thank the sponsoring organizations: Microsoft Research Asia, the Department of Systems Engineering and Engineering Management at the Chinese University of Hong Kong, and LexisNexis for their financial support, the Department of Computer Science and Technology, Tsinghua University for local arrangements, the Chinese NewsML Community for website design and administration, Ling Huang for the logistics, Weiwei Sun for the conference webpage management, EONSOLUTION for the conference management, and Springer for the postconference

LNCS publication. We believe that this conference set a very high standard for a regionally oriented conference, especially in Asia, and we hope that it continues as a tradition in the upcoming years.


Sung Hyon Myaeng and Ming Zhou (PC Co-chairs)
Kam-Fai Wong and Hong-Jiang Zhang (Conference Co-chairs)

# Organization

## General Conference Co-chairs

Kam-Fai Wong, Chinese University of Hong Kong, China
Hong-Jiang Zhang, Microsoft Research Asia, China

## Program Co-chairs

Sung Hyon Myaeng, Information and Communications University (ICU),
    South Korea
Ming Zhou, Microsoft Research Asia, China

## Organization Committee Co-chairs

Jianfeng Gao, Microsoft Research Asia, China
Shao-Ping Ma, Tsinghua University, China

## Special Session and Tutorials Chair

Wei-Ying Ma, Microsoft Research Asia, China

## Interactive Posters/Demo Chair

Gary Geunbae Lee, POSTECH, South Korea

## Steering Committee

Jun Adachi, National Institute of Informatics, Japan
Hsin-Hsi Chen, National Taiwan University, Taiwan
Lee-Feng Chien, Academia Sinica, Taiwan
Tetsuya Ishikawa, University of Tsukuba, Japan
Gary Geunbae Lee, POSTECH, South Korea
Mun-Kew Leong, Institute for Infocomm Research, Singapore
Helen Meng, Chinese University of Hong Kong, China
Sung-Hyon Myaeng, Information and Communications University, South Korea

Hwee Tou Ng, National University of Singapore, Singapore
Kam-Fai Wong, Chinese University of Hong Kong, China

## Organizing Committee

Lee-Feng Chien, Academia Sinica, Taiwan (Publicity, Asia)
Susan Dumais, Microsoft, USA (Publicity, North America)
Jianfeng Gao, Microsoft Research Asia, China (Publication, Co-chair)
Mun-Kew Leong, Institute for Infocomm Research, Singapore (Finance)
Shao-Ping Ma, Tsinghua University, China (Local Organization, Co-chair)
Ricardo Baeza-Yates, University of Chile, Chile (Publicity, South America)
Shucai Shi, BITI, China (Local Organization)
Dawei Song, DSTC, Australia (Publicity, Australia)
Ulich Thiel, IPSI, Germany (Publicity, Europe)
Chuanfa Yuan, Tsinghua University, China (Local Organization)

# Program Committee

Peter Anick, Yahoo, USA
Hsin-Hsi Chen, National Taiwan University, Taiwan
Aitao Chen, University of California, Berkeley, USA
Lee-Feng Chien, Academia Sinica, Taiwan
Fabio Crestani, University of Strathclyde, UK
Edward A. Fox, Virginia Tech, USA
Jianfeng Gao, Microsoft Research Asia, China
Hani Abu-Salem, DePaul University, USA
Tetsuya Ishikawa, University of Tsukuba, Japan
Christopher Khoo, Nanyang Technological University, Singapore
Jung Hee Kim, North Carolina A&T University, USA
Minkoo Kim, Ajou University, South Korea
Munchurl Kim, Information and Communication University, South Korea
Kazuaki Kishida, Surugadai University, Japan
Kui-Lam Kwok, Queens College, City University of New York, USA
Wai Lam, Chinese University of Hong Kong, China
Gary Geunbae Lee, POSTECH, South Korea
Mun-Kew Leong, Institute for Infocomm Research, Singapore
Gena-Anne Levow, University of Chicago, USA
Hang Li, Microsoft Research Asia, China
Robert Luk, Hong Kong Polytechnic University, China
Gay Marchionini, University of North Carolina, Chapel Hill, USA
Helen Meng, Chinese University of Hong Kong, China
Hiroshi Nakagawa, University of Tokyo, Japan
Hwee Tou Ng, National University of Singapore, Singapore
Jian-Yun Nie, University of Montreal, Canada
Jon Patrick, University of Sydney, Australia
Ricardo Baeza-Yates, University of Chile, Chile
Hae-Chang Rim, Korea University, South Korea
Tetsuya Sakai, Toshiba Corporate R&D Center, Japan
Padmini Srinivasan, University of Iowa, USA
Tomek Strzalkowski, State University of New York, Albany, USA
Maosong Sun, Tsinghua University, China
Ulrich Thiel, Fraunhofer IPSI, Germany
Takenobu Tokunaga, Tokyo Institute of Technology, Japan
Hsin-Min Wang, Academia Sinica, Taiwan
Ross Wilkinson, CSIRO, Australia
Lide Wu, Fudan University, China
Jinxi Xu, BBN Technologies, USA
ChengXiang Zhai, University of Illinois, Urbana Champaign, USA
Min Zhang, Tsinghua University, China

## Reviewers

| | | |
|---|---|---|
| Peter Anick | Kui-Lam Kwok | Hae-Chang Rim |
| Yunbo Cao | Wai Lam | Tetsuya Sakai |
| Yee Seng Chan | Gary Geunbae Lee | Tomek Strzalkowski |
| Hsin-Hsi Chen | Mun-Kew Leong | Maosong Sun |
| Zheng Chen | Gena-Anne Levow | Ulrich Thiel |
| Aitao Chen | Hang Li | Takenobu Tokunaga |
| Tee Kiah Chia | Mu Li | Hsin-Min Wang |
| Lee-Feng Chien | Hongqiao Li | Haifeng Wang |
| Fabio Crestani | Chin-Yew Lin | Ross Wilkinson |
| Edward A. Fox | Robert Luk | Kam-Fai Wong |
| Jianfeng Gao | Wei-Ying Ma | Lide Wu |
| Hani Abu-Salem | Gay Marchionini | Jinxi Xu |
| Xuanjing Huang | Helen Meng | Peng Yu |
| Tetsuya Ishikawa | Sung Hyon Myaeng | Chunfa Yuan |
| Christopher Khoo | Hiroshi Nakagawa | ChengXiang Zhai |
| Jung Hee Kim | Hwee Tou Ng | Hong-Jiang Zhang |
| Minkoo Kim | Jian-Yun Nie | Min Zhang |
| Munchurl Kim | Jon Patrick | Ming Zhou |
| Kazuaki Kishida | Ricardo Baeza-Yates | Jian-lai Zhou |

# Table of Contents

## Information Organization

## Automatic Summarization

## Alignment/Paraphrasing in IR

## Web Search

## Linguistic Issues in IR

## Document/Query Models

## Enabling Technology

## Mobile Applications