# Preface

CICLing-2004 is the fifth Annual Conference on Intelligent Text Processing and Computational Linguistics, see www.CICLing.org. CICLing conferences are intended to provide a balanced view of the cutting edge developments in both theoretical foundations of computational linguistics and practice of natural language text processing with its numerous applications. A feature of CICLing conferences is their wide scope that covers nearly all areas of computational linguistics and all aspects of natural language processing applications. These conferences are a forum for dialogue between the specialists working in the two areas.

This year we were honored by the presence of our invited speakers *Martin Kay* of Stanford University, *Philip Resnik* of the University of Maryland, *Ricardo Baeza-Yates*, of the University of Chile, and *Nick Campbell* of the ATR Spoken Language Translation Research Laboratories. They delivered excellent extended lectures and organized vivid discussions.

Of 129 submissions received (74 full papers and 44 short papers), after careful international reviewing 74 papers have been selected for presentation (40 full papers and 35 short papers), by 176 authors from 21 countries: Korea (37), Spain (34), Japan (22), Mexico (15), China (11), Germany (10), Ireland (10), UK (10), Singapore (6), Canada (3), Czech Rep. (3), France (3), Brazil (2), Sweden (2), Taiwan (2), Turkey (2), USA (2), Chile (1), Romania (1), Thailand (1), The Netherlands (1); the figures in parentheses stand for the number of authors from the corresponding country.

In addition to high scientific level, one of the success factors of CICLing conferences is their excellent cultural program. CICLing-2004 was held in Korea, the beautiful and wonderful Country of the Morning Calm, as Korean people call their land. The participants enjoyed three full-day excursions to the most important natural and historical attractions around Seoul city; see photos at www.CICLing.org. Full-day excursions allow for friendly personal interaction between participants and give them a chance to make friends with the most famous experts in the field, who are not easily accessible at larger conferences.

A conference is the result of the work of many people. First of all I would like to thank the members of the Program Committee for the time and effort they devoted to the reviewing of the submitted articles and to the selection process. Especially helpful were Manuel Vilares, John Tait, Alma Kharrat, Karin Verspoor, Viktor Pekar, and many others—a complete list would be too long.

Obviously I thank the authors for their patience in the preparation of the papers, not to mention the very development of their scientific results that form this book. I also express my most cordial thanks to the members of the local Organizing Committee for their considerable contribution to making this conference become a reality. Last but not least, I thank our host—the ITRI of the Chung-Ang University. would like to also thank RITOS-2 of CYTED for their support of the CICLing conferences.

December 2003                                                            Alexander Gelbukh

## Conference Chair

Alexander Gelbukh (CIC-IPN, Mexico / Chung-Ang U., Korea)


## Program Committee

1. Boitet, Christian (CLIPS-IMAG, France),
2. Bolshakov, Igor (CIC-IPN, Mexico),
3. Bontcheva, Kalina (U. Sheffield, UK),
4. Calzolari, Nicoletta (ILC-CNR, Italy),
5. Carroll, John, (U. Sussex, UK)
6. Cristea, Dan (U. Iasi, Romania),
7. Gelbukh, Alexander (**Chair**, CIC-IPN, Mexico / Chung-Ang U., Korea)
8. Hallett, Cătălina (U. Brighton, UK),
9. Han, SanYong (Chung-Ang U., Korea)
10. Harada, Yasunari (Waseda U, Japan),
11. Hasida, Kôiti (Electrotechnical Laboratory-AIST, Japan),
12. Hirst, Graeme (U. Toronto, Canada),
13. Hovy, Eduard (ISI of U. Southern Carolina, USA),
14. Johnson, Frances (Manchester Metropolitan U., UK),
15. Kharrat, Alma (Microsoft Research, USA),
16. Kilgarriff, Adam (U. Brighton, UK),
17. Kittredge, Richard (CoGenTex Inc., USA / Canada),
18. Kübler, Sandra (U. Tübingen, Germany),
19. López López, Aurelio (INAOE, Mexico),
20. Loukanova, Roussanka (Indiana U, USA / Bulgaria),
21. Lüdeling, Anke (U. Stuttgart, Germany),
22. Maegard, Bente (Centre for Language Technology, Denmark),
23. Martín-Vide, Carlos (U. Rovira i Virgili, Spain),
24. Mel'čuk, Igor (U. Montreal, Canada),
25. Metais, Elisabeth (U. Versailles, France),
26. Mihalcea, Rada (U. North Texas,  USA),
27. Mitkov, Ruslan (U. Wolverhampton, UK),
28. Murata, Masaki (KARC-CRL, Japan),
29. Narin'yani, Alexander (Russian Institute of Artificial Intelligence, Russia),
30. Nirenburg, Sergei (New Mexico U, USA),
31. Palomar, Manuel (U. Alicante, USA / Spain),
32. Pedersen, Ted (U. Minnesota Duluth, USA),
33. Pekar, Viktor (U. Wolverhampton, UK)
34. Pineda Cortes, Luis Alberto (UNAM, Mexico),
35. Piperidis, Stelios (Institute for Language and Speech Processing, Greece),
36. Pustejovsky, James (Brandeis U., USA)
37. Ren, Fuji (U. Tokushima, Japan),
38. Riloff, Ellen (U. Utah, USA)
39. Sag, Ivan (Standford U, USA),
40. Sharoff, Serge (U. Leeds, UK),

41. Sidorov, Grigori (CIC-IPN, Mexico),
42. Sun Maosong (Tsinghua U, China),
43. Tait, John (U. Sunderland, UK),
44. Trujillo, Arturo (UMIST, UK),
45. T'sou Ka-yin, Benjamin (City U. Hong Kong, Hong Kong),
46. Van Guilder, Linda (MITRE Corp., USA),
47. Verspoor, Karin (Intelligenesis Corp., USA / The Netherlands),
48. Vilares Ferro, Manuel (U. La Coruña, Spain),
49. Wilks, Yorick (U. Sheffield, UK).

## Additional Reviewers

1. Babych, Bogdan (U. Leeds, Centre for Translation, UK)
2. Campbell, Nick (ATR Human Information Science Labs, Japan)
3. Liang Shao-Fen (U. Sunderland; Taiwan)
4. Llopis, Fernando (U. Alicante, Spain)
5. Martínez-Barco, Patricio (U. Alicante, Spain)
6. Montoyo, Andrés (U. Alicante, Spain)
7. Oakes, Michael (U. Sunderland; Taiwan)
8. Saiz Noeda, Maximiliano (U. Alicante, Spain)
9. Stokoe, Christopher (U. Sunderland, UK)
10. Vicedo, José L. (U. Alicante, Spain)

## Organizing Committee

1. Han, SanYong (**chair**)
2. Gelbukh, Alexander (**co-chair**)
3. Chang, Tae Gyu
4. Shim, Duk Sun
5. Kim, Jun Sung
6. Park, Ho Hyun
7. Kim, Hyung Suk
8. Moon, Suk Whan
9. Choi, Ok Kyung
10. Kang, Nam Oh
11. Shin, Kwang Chul

## Organization. Website and Contact

The conference was organized by the Electronic Commerce and Internet Applications Laboratory of the Chung-Ang University, Seoul, Korea, and the Natural Language Laboratory of the Center for Computing Research of the National Polytechnic

Institute, Mexico City, Mexico. The conference was hosted by the ITRI of the Chung-Ang University.

The website of the CICLing conferences is www.CICLing.org (mirrored at www.cic.ipn.mx/cicling). Contact: CICLing.org; see also www.gelbukh.com.

# Table of Contents

## Computational Linguistics

## Bilingual Resources

## Machine Translation

## Natural Language Generation

## Human-Computer Interaction Applications

## Speech Recognition and Synthesis

## Intelligent Text Processing

## Indexing

# Information Extraction

# Text Categorization

# Document Clustering

# Summarization