# Theory and Applications of Natural Language Processing

Series Editors:
Graeme Hirst (Textbooks)
Eduard Hovy (Edited volumes)
Mark Johnson (Monographs)

## Aims and Scope

The field of Natural Language Processing (NLP) has expanded explosively over the past decade: growing bodies of available data, novel fields of applications, emerging areas and new connections to neighboring fields have all led to increasing output and to diversification of research.

"Theory and Applications of Natural Language Processing" is a series of volumes dedicated to selected topics in NLP and Language Technology. It focuses on the most recent advances in all areas of the computational modeling and processing of speech and text across languages and domains. Due to the rapid pace of development, the diversity of approaches and application scenarios are scattered in an ever-growing mass of conference proceedings, making entry into the field difficult for both students and potential users. Volumes in the series facilitate this first step and can be used as a teaching aid, advanced-level information resource or a point of reference.

The series encourages the submission of research monographs, contributed volumes and surveys, lecture notes and textbooks covering research frontiers on all relevant topics, offering a platform for the rapid publication of cutting-edge research as well as for comprehensive monographs that cover the full range of research on specific problem areas.

The topics include applications of NLP techniques to gain insights into the use and functioning of language, as well as the use of language technology in applications that enable communication, knowledge management and discovery such as natural language generation, information retrieval, question-answering, machine translation, localization and related fields.

The books are available in printed and electronic (e-book) form:

* Downloadable on your PC, e-reader or iPad
* Enhanced by Electronic Supplementary Material, such as algorithms, demonstrations, software, images and videos
* Available online within an extensive network of academic and corporate R&D libraries worldwide
* Never out of print thanks to innovative print-on-demand services
* Competitively priced print editions for eBook customers thanks to MyCopy service http://www.springer.com/librarians/e-content/mycopy

Antal van den Bosch • Gosse Bouma
Editors

# Interactive Multi-modal Question-Answering

*Editors*

Antal van den Bosch
Tilburg center for Cognition and Communication
Tilburg University
School of Humanities
5000 LE Tilburg
The Netherlands
Antal.vdnBosch@uvt.nl

Gosse Bouma
Information Science
University of Groningen
9700 AS Groningen
The Netherlands
g.bouma@rug.nl

# Preface

To start off with a question, what kind of book is this? This book is, in many aspects, a collaborative effort. Its contents evolve from the idea of combining the latest research in fields such as natural language processing, dialogue systems, human-machine interaction, and information extraction, around a single common theme. What is more, the book documents a successful collaborative effort to do just this: the combination of all these fields in one single research area, working towards a common goal—an interactive, multimodal spoken question-answering dialogue system, with which a human user could discuss a particular domain. This concerted effort was known as the "Interactive multimodal Information eXtraction" (IMIX) Programme, a Dutch national research programme funded by the Netherlands Organisation for Scientific Research (NWO), that ran between 2004 and 2009.

During the peak years of IMIX, the researchers in this project together formed a major part of the entire language and speech technology community in the Netherlands. The project also continued to have an influence in the work that followed, for example in the post-doc research tracks of the people who wrote their PhD thesis as part of IMIX. The majority of researchers involved in IMIX can be traced back as authors of the chapters in this book and, although the selection of chapters remains incomplete with respect to all the work performed in IMIX, we are very grateful to the many former IMIX colleagues who were willing to write retrospective overviews of their work and performance in IMIX.

The IMIX project was surrounded by a supportive network of individuals, each of whom deserves a special mention: Alice Dijkstra, Christine Erb, and Brigit van der Pas from NWO were instrumental throughout the project in organising regular IMIX meetings and, more generally, in creating an atmosphere in which everyone was encouraged to maximise performance and collaborate with as few constraints and restrictions as possible. Alice Dijkstra, who herself has a background in computational linguistics, had created a sound foundation by inviting a team of academic strategists, headed by Lou Boves, to write up the most challenging mix of ideas that, within reason, could be expected to be illustrated by the Dutch language and speech technology community, on a limited budget and within the time span of a PhD project. After an intensive assessment proposal period, the IMIX programme committee settled on the project proposals. The end results of most of these projects are covered in this book. The programme committeee appointed a project coordinator, a post occupied from the start by Els den Os, to oversee the creation of an actual running, talking, question-answering piece of demonstrator software. This book also highlights this complicated endeavour, which to some extent resembled putting together a working car from a washing machine, two lawnmowers, and a microwave oven (all ultra modern).

The project was fortunate to have the guidance of an international scientific advisory board: Eduard Hovy, Jon Oberlander, Norbert Reithinger, and Steve Young. From start to finish, their advice gave the project all the right nudges, in the right directions. Their final report now appears as the epilogue of this book. On behalf of the whole IMIX team, we would like to express our gratitude to our

international advisers and to everyone who helped make this challenging project possible. For assisting us in the final editorial work, we thank Olga Chiarcos at Springer, and the team of copyeditors from CumLingua.

This book is the final publication from the IMIX project. Its aims to go beyond IMIX; it is intended as a unique look at a mix of state-of-the-art methods, rarely found together in one book (despite the fact that we are all involved in language and speech technology—such is the level of specialisation of most of that we do). But in addition to all the developments and future work born from the ideas orig in IMIX, it gives us tremendous pleasure to have had the privilege to edit this book and provide a final answer to the first question: this book is IMIX.

Tilburg and Groningen                          Antal van den Bosch and Gosse Bouma

# Contents

## Part II  Interaction Management

## Part III  Fusing Text, Speech, and Images

**Part IV   Text Analysis for Question Answering**