

Lecture Notes in Artificial Intelligence 6001

Edited by R. Goebel, J. Siekmann, and W. Wahlster

Subseries of Lecture Notes in Computer Science

Thiago Alexandre Salgueiro Pardo
António Branco Aldebaro Klautau
Renata Vieira Vera Lúcia Strube de Lima (Eds.)

Computational Processing of the Portuguese Language

9th International Conference, PROPOR 2010
Porto Alegre, RS, Brazil, April 27-30, 2010
Proceedings

Series Editors

Randy Goebel, University of Alberta, Edmonton, Canada

Jörg Siekmann, University of Saarland, Saarbrücken, Germany

Wolfgang Wahlster, DFKI and University of Saarland, Saarbrücken, Germany

Volume Editors

Thiago Alexandre Salgueiro Pardo

Universidade de São Paulo, São Carlos / SP, Brazil

E-mail: taspardo@icmc.usp.br

António Branco

Universidade de Lisboa, Portugal

E-mail: antonio.branco@di.fc.ul.pt

Aldebaro Klautau

Universidade Federal do Pará, Belém / PA, Brazil

E-mail: aldebaro@ufpa.br

Renata Vieira

PUCRS, Porto Alegre / RS, Brazil

E-mail: renata.vieira@pucrs.br

Vera Lúcia Strube de Lima

PUCRS, Porto Alegre / RS, Brazil

E-mail: vera.strube@pucrs.br

Library of Congress Control Number: 2010923533

CR Subject Classification (1998): I.2, H.3, H.4, I.4, I.5, H.2.8

LNCS Sublibrary: SL 7 – Artificial Intelligence

ISSN 0302-9743

ISBN-10 3-642-12319-8 Springer Berlin Heidelberg New York

ISBN-13 978-3-642-12319-1 Springer Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

springer.com

© Springer-Verlag Berlin Heidelberg 2010
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India
Printed on acid-free paper 06/3180

Preface

The International Conference on Computational Processing of Portuguese—PROPOR—is the main event in the area of natural language processing that is focused on Portuguese and the theoretical and technological issues related to this language. It welcomes contributions for both written and spoken language processing.

The event is hosted in Brazil and in Portugal. The meetings have been held in Lisbon/Portugal (1993), Curitiba/Brazil (1996), Porto Alegre/Brazil (1998), Évora/Portugal (1999), Atibaia/Brazil (2000), Faro/Portugal (2003), Itatiaia/Brazil (2006) and Aveiro/Portugal (2008).

This meeting has been a highly productive forum for the progress of this area and to foster the cooperation among the researchers working on the automated processing of the Portuguese language. PROPOR brings together research groups, promoting the development of methodologies, resources and projects that can be shared among all researchers and practitioners in the field.

The ninth edition of this event was held in Porto Alegre, Brazil, at *Pontifícia Universidade Católica do Rio Grande do Sul* (PUCRS). It had two main tracks: one for language processing and another one for speech processing. This event hosted a special Demonstration Session and the first edition of the PhD and MSc Dissertation Contest, which aimed at recognizing the best academic work on processing of the Portuguese language in the last few years. This edition of the event featured tutorials on statistical machine translation and on speech recognition, as well as invited talks by renowned researchers of natural language processing.

A total of 48 submissions were received, 37 for the language track and 11 for the speech track, by authors from 10 countries: Brazil, China, Denmark, UK, Germany, Italy, Poland, Portugal, Spain and USA. Each submission was evaluated by at least three members from a multidisciplinary and international scientific committee.

This volume gathers a selection of the 21 best papers accepted to be presented at this meeting, of which 13 are full papers, corresponding to an acceptance rate of 27%. These papers cover the areas concerning applications for information handling and text processing, language processing, language resources, and speech recognition and synthesis.

We would like to express our thanks to everyone involved in the organization of the event, to the scientific committee members for their excellent work, to the researchers who kindly accepted to contribute to the event by delivering tutorials and invited talks, and to the institutions, organizations and funding agencies which allowed the realization of this event, namely, PUCRS, SBC (the Brazilian Computer Society), CEPLN (the SBC Special Interest Group on Natural Language Processing), CNPq (*Conselho Nacional de Desenvolvimento Científico e Tecnológico* – a Brazilian funding agency), NAACL (The North American Chapter of the Association for Computational Linguistics), ISCA (International Speech Communication Association), SIG-IL (the ISCA Special Interest Group on Iberian Languages) and CLARIN

(Common Language Resources and Technology Infrastructure – a large-scale pan-European collaborative effort to create, coordinate and make language resources and technology available and readily usable).

April 2010

Thiago Alexandre Salgueiro Pardo
António Branco
Aldebaro Klautau
Renata Vieira
Vera Lúcia Strube de Lima

Organization

General Chair

Vera Lúcia Strube de Lima Pontifícia Universidade Católica do Rio Grande do Sul, Brazil

Program Chairs

António Branco Universidade de Lisboa, Portugal - Language
Aldebaro Klautau Universidade Federal do Pará, Brazil - Speech

Tutorial Chair

Maria das Graças Volpe Nunes Universidade de São Paulo, Brazil

Editorial Chair

Thiago Alexandre Salgueiro Pardo Universidade de São Paulo, Brazil

PhD and MSc Dissertation Contest Chair

David de Matos Instituto de Engenharia de Sistemas e Computadores, Portugal

Demo Session Chair

António Teixeira Universidade de Aveiro, Portugal

Local Organizing Chair

Renata Vieira Pontifícia Universidade Católica do Rio Grande do Sul, Brazil

Program Committee

Alexandre Agustini

Pontifícia Universidade Católica do Rio Grande do Sul, Brazil

Aline Villavicencio

Universidade Federal do Rio Grande do Sul, Brazil

Amália Mendes

Universidade de Lisboa, Portugal

Andre Adami

Universidade de Caxias do Sul, Brazil

António Ribeiro

European Commission, Joint Research Centre, Italy

António Serralheiro

Instituto de Engenharia de Sistemas e Computadores /

António Teixeira

Instituto Superior Técnico, Portugal

Antonio Bonafonte

Universidade de Aveiro, Portugal

Augusto Silva

Universitat Politècnica de Catalunya, Spain

Bento Dias da Silva

Universidade Católica Portuguesa, Portugal

Berthold Crysmann

Universidade Estadual Paulista, Brazil

Carlos Prolo

Universität Bonn / Universität des Saarlandes,

Caroline Hagège

Germany

Celso Kaestner

Pontifícia Universidade Católica do Rio Grande do

Climent Nadeu

Sul, Brazil

Cristina Martins

Xerox Research Centre, France

Daniela Braga

Universidade Tecnológica Federal do Paraná, Brazil

Dante Barone

Universitat Politècnica de Catalunya, Spain

David de Matos

Universidade de Coimbra, Portugal

Diamantino Freitas

Microsoft Language Development Center, Portugal

Edmilson Morais

Universidade Federal do Rio Grande do Sul, Brazil

Eric Laporte

Instituto de Engenharia de Sistemas e Computadores,

Fábio Violaro

Portugal

Fernando Resende

Universidade do Porto, Portugal

Florence Amardeilh

Vocalize, Brazil

Gabriel Pereira Lopes

Université Paris-Est Marne-la-Vallée, France

Gael Harry Dias

Universidade Estadual de Campinas, Brazil

Horacio Saggion

Universidade Federal do Rio de Janeiro, Brazil

Irene Rodrigues

Mondeca, France

Isabel Trancoso

Universidade Nova de Lisboa, Portugal

Ivandro Paraboni

Universidade da Beira Interior, Portugal

Ivandro Sanches

University of Sheffield, UK

Jason Baldridge

Universidade de Évora, Portugal

Jean-Luc Minel

Instituto de Engenharia de Sistemas e Computadores /

João Balsa

Instituto Superior Técnico, Portugal

João Paulo Neto

Universidade de São Paulo, Brazil

Jorge Baptista

Centro Universitário da FEI, Brazil

Julia Hirschberg

University of Texas, USA

Université Paris Ouest - Nanterre La Défense, France

Universidade de Lisboa, Portugal

Instituto de Engenharia de Sistemas e Computadores,

Portugal

Universidade do Algarve, Portugal

Columbia University, USA

Lúcia Rino	Universidade Federal de São Carlos, Brazil
Luísa Coheur	Instituto de Engenharia de Sistemas e Computadores, Portugal
Luiz Pizzato	University of Sydney, Australia
Marcelo Finger	Universidade de São Paulo, Brazil
Marco Gonzalez	Pontifícia Universidade Católica do Rio Grande do Sul, Brazil
Maria Viana	Universidade de Lisboa, Portugal
Maria das Graças Volpe	Universidade de São Paulo, Brazil
Nunes	
Maria Francisca Xavier	Universidade Nova de Lisboa, Portugal
Maximiliano Saiz Noeda	Universidad de Alicante, Spain
Miguel Filgueiras	Universidade do Porto, Portugal
Nestor Yoma	Universidad de Chile, Chile
Nuno Mamede	Instituto de Engenharia de Sistemas e Computadores / Instituto Superior Técnico, Portugal
Nuno Cavalheiro Marques	Universidade Nova de Lisboa, Portugal
Pablo Gamallo	Universidad de Santiago de Compostela, Spain
Paulo Quaresma	Universidade de Évora, Portugal
Plínio Barbosa	Universidade Estadual de Campinas, Brazil
Ranniere Maia	Toshiba Research Europe Limited, UK
Renata Vieira	Pontifícia Universidade Católica do Rio Grande do Sul, Brazil
Sandra Aluisio	Universidade de São Paulo, Brazil
Stanley Loh	Universidade Católica de Pelotas, Brazil
Steven Bird	University of Melbourne, Australia
Thiago Alexandre Salgueiro	Universidade de São Paulo, Brazil
Pardo	
Tracy King	Microsoft, USA
Violeta Quental	Pontifícia Universidade Católica do Rio de Janeiro, Brazil
Viviane Orengo	Universidade Federal do Rio Grande do Sul, Brazil

Additional Reviewers

Aldebaro Klautau	Universidade Federal do Pará, Brazil
Belinda Maia	Universidade do Porto, Portugal
Fábio Kepler	Universidade de São Paulo, Brazil
Guillaume Jacquet	Xerox Research Centre, France
Jussara Vieira	Vocalize, Brazil
Marta Costa-jussà	Barcelona Media, Spain
Milagros Fernández	Universidade da Coruña, Spain
Gavilanes	
Simon Zwarts	Macquarie University, Australia

Steering Committee

Maria das Graças Volpe Nunes	Universidade de São Paulo, Brazil
Vera Lucia Strube de Lima	Pontifícia Universidade Católica do Rio Grande do Sul, Brazil
Isabel Trancoso	Instituto de Engenharia de Sistemas e Computadores / Instituto Superior Técnico, Portugal
Violeta Quental	Pontifícia Universidade Católica do Rio de Janeiro, Brazil
António Teixeira	Universidade de Aveiro, Portugal

Table of Contents

Applications: Information Handling

Improving IdSay: A Characterization of Strengths and Weaknesses in Question Answering Systems for Portuguese	1
<i>Gracinda Carvalho, David Martins de Matos, and Vitor Rocio</i>	
Assessing the Impact of Stemming Accuracy on Information Retrieval	11
<i>Felipe N. Flores, Viviane P. Moreira, and Carlos A. Heuser</i>	
Exploiting Multilingual Grammars and Machine Learning Techniques to Build an Event Extraction System for Portuguese	21
<i>Vanni Zavarella, Hristo Tanev, Jens Linge, Jakub Piskorski, Martin Atkinson, and Ralf Steinberger</i>	
Formalizing CST-Based Content Selection Operations	25
<i>Maria Lucia Castro Jorge and Thiago Alexandre Salgueiro Pardo</i>	

Applications: Text Processing

Translating from Complex to Simplified Sentences	30
<i>Lucia Specia</i>	
Challenging Choices for Text Simplification	40
<i>Caroline Gasperin, Erick Maziero, and Sandra M. Aluísio</i>	
Comparing Sentence-Level Features for Authorship Analysis in Portuguese	51
<i>Rui Sousa-Silva, Luís Sarmento, Tim Grant, Eugénio Oliveira, and Belinda Maia</i>	

Language Processing

A Machine Learning Approach to Portuguese Clause Identification	55
<i>Eraldo R. Fernandes, Cícero N. dos Santos, and Ruy L. Milidiú</i>	
A Hybrid Approach for Multiword Expression Identification	65
<i>Carlos Ramisch, Helena de Medeiros Caseli, Aline Villavicencio, André Machado, and Maria José Finatto</i>	
Out-of-the-Box Robust Parsing of Portuguese	75
<i>João Silva, António Branco, Sérgio Castro, and Ruben Reis</i>	
LXGram: A Deep Linguistic Processing Grammar for Portuguese	86
<i>Francisco Costa and António Branco</i>	

Language Resources

InferenceNet.Br: Expression of Inferentialist Semantic Content of the Portuguese Language	90
<i>Vladia Pinheiro, Tarcisio Pequeno, Vasco Furtado, and Wellington Franco</i>	
Comparing Verb Synonym Resources for Portuguese	100
<i>Jorge Teixeira, Luís Sarmento, and Eugénio Oliveira</i>	
Auxiliary Verbs and Verbal Chains in European Portuguese	110
<i>Jorge Baptista, Nuno Mamede, and Fernando Gomes</i>	
P-AWL: Academic Word List for Portuguese	120
<i>Jorge Baptista, Neuza Costa, Joaquim Guerra, Marcos Zampieri, Maria Cabral, and Nuno Mamede</i>	

Speech Recognition

Automatic Phone Clustering Based on Confusion Matrices	124
<i>Carla Lopes, Arlindo Veiga, and Fernando Perdigão</i>	
An Open-Source Speech Recognizer for Brazilian Portuguese with a Windows Programming Interface	128
<i>Patrick Silva, Pedro Batista, Nelson Neto, and Aldebaro Klautau</i>	
A Baseline System for Continuous Speech Recognition of Brazilian Portuguese Using the West Point Brazilian Portuguese Speech Corpus	132
<i>Fabiano Weimar dos Santos, Dante Augusto Couto Barone, and André Gustavo Adami</i>	

Speech Synthesis

Voice Quality of European Portuguese Emotional Speech	142
<i>Ana Nunes, Rosa Lídia Coimbra, and António Teixeira</i>	
Prosodic Prediction in Brazilian Portuguese: A Contribution to Speech Synthesis	152
<i>Cirineu Cecote Stein</i>	
The Role of Morphology in Generating High-Quality Pronunciation Lexica for Regional Variants of Portuguese	162
<i>Simone Ashby and José Pedro Ferreira</i>	

Author Index	167
---------------------------	------------