

Lecture Notes in Artificial Intelligence 7243

Subseries of Lecture Notes in Computer Science

LNAI Series Editors

Randy Goebel

University of Alberta, Edmonton, Canada

Yuzuru Tanaka

Hokkaido University, Sapporo, Japan

Wolfgang Wahlster

DFKI and Saarland University, Saarbrücken, Germany

LNAI Founding Series Editor

Joerg Siekmann

DFKI and Saarland University, Saarbrücken, Germany

Helena Caseli Aline Villavicencio
António Teixeira Fernando Perdigão (Eds.)

Computational Processing of the Portuguese Language

10th International Conference, PROPOR 2012
Coimbra, Portugal, April 17-20, 2012
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

Helena Caseli
UFSCAR, São Carlos, SP, Brazil
E-mail: helenacaseli@dc.ufscar.br

Aline Villavicencio
UFRGS, Porto Alegre, RS, Brazil
E-mail: alinev@gmail.com

António Teixeira
Universidade de Aveiro, Aveiro, Portugal
E-mail: ajst@ua.pt

Fernando Perdigão
Universidade de Coimbra, Coimbra, Portugal
E-mail: fp@deec.uc.pt

ISSN 0302-9743 e-ISSN 1611-3349
ISBN 978-3-642-28884-5 e-ISBN 978-3-642-28885-2
DOI 10.1007/978-3-642-28885-2
Springer Heidelberg Dordrecht London New York

Library of Congress Control Number: 2012933122

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

LNCS Sublibrary: SL 7 – Artificial Intelligence

© Springer-Verlag Berlin Heidelberg 2012

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.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

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), Aveiro/Portugal (2008), and Porto Alegre/Brazil (2010).

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 tenth edition of this event was held at the University of Coimbra, Coimbra, Portugal. It had two main tracks: one for language processing and another for speech processing. This event hosted a special Demonstration Session and a satellite event named “Págico,” consisting in an evaluation contest for non-trivial information seeking in Portuguese using Wikipedia as target. This edition of PROPOR featured two invited talks by internationally renowned researchers as well as tutorials on symbolic and statistical approaches to natural language processing and analysis and visual feedback of the singing voice.

A total of 86 submissions were received, 61 for the language track and 25 for the speech track, by authors in worldwide institutions from countries like Brazil, China, Germany, Portugal, and Spain. Each submission was evaluated by at least three members from a multidisciplinary and international scientific committee.

This volume gathers a selection of the 47 best papers accepted to be presented at this meeting, of which 24 are full papers, corresponding to an acceptance rate of 27%. These papers cover the areas related to automatic acquisition of information, linguistic description and processing, language resources, language applications and speech production, speech processing and applications.

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, University of Coimbra, IT (Instituto de Telecomunicações), FCT (Portuguese National Founding Agency), ISCA (International Speech Communication Association), SIG-IL (the ISCA Special Group on Iberian Languages), CEPLN (the Special Interest Group on Natural Language Processing of the Brazilian Computer Society), and ACL (Association for Computational Linguistics).

April 2012

Helena de Medeiros Caseli

Aline Villavicencio

António Teixeira

Fernando Perdigão

Organization

General Chair

Fernando Perdigão	Universidade de Coimbra/IT, Portugal
-------------------	--------------------------------------

Program Chairs

Aline Villavicencio	UFRGS, Brazil – Language
António Teixeira	Universidade de Aveiro/IEETA, Portugal – Speech

Editorial Chair

Helena de Medeiros Caseli	UFSCAR, Brazil
---------------------------	----------------

Demo Session Chair

Alberto Abad	L2F INESC-ID, Portugal
--------------	------------------------

PhD and MSc/MA Dissertation Contest Chair

Jorge Baptista	Universidade do Algarve, Portugal
----------------	-----------------------------------

Local Organizing Committee

Luís Sá	UC/IT, Portugal
Sara Candeias	IT, Portugal
Ana R. Luís	UC/CELGA, Portugal
Carla Lopes	IPLEI/IT, Portugal
Hugo Gonçalo Oliveira	UC/CISUC, Portugal

Program Committee

Alberto Abad	INESC-ID, Portugal
Alberto Simões	UM, Portugal
Aldebaro Klautau	UFPA, Brazil
Alexandre Agustini	PUC-RS, Brazil
Aline Villavicencio	UFRGS, Brazil
Amália Andrade	UL, Portugal
Amália Mendes	UL, Portugal

Ana Bocorny	PUC-RS, Brazil
Ana Luís	UC, Portugal
Andre Adami	UCS, Brazil
Andreia Rauber	UCPEL, Brazil
Antonio Bonafonte	UPC, Spain
António Branco	UL, Portugal
Antonio Rubio	UG, Spain
António Serralheiro	INESC-ID, Portugal
António Teixeira	Universidade de Aveiro, Portugal
Ariadne Carvalho	Unicamp, Brazil
Ariani Di Felippo	UFSCAR, Brazil
Belinda Maia	UP, Portugal
Bento da Silva	UNESP, Brazil
Berthold Crysmann	CNRS Paris-Diderot, France
Carlos Prolo	PUC-RS, Brazil
Carlos Teixeira	UL, Portugal
Carmen García Mateo	UV, Spain
Caroline Gasperin	TouchType, UK
Caroline Hagège	Xerox Research Centre, France
Catarina Oliveira	Universidade de Aveiro, Portugal
Ciro Martins	Universidade de Aveiro, Portugal
Cristiane Killian	UFRGS, Brazil
Daniela Braga	Microsoft, China
Dante Barone	UFRGS, Brazil
Diana Santos	University of Oslo, Norway
Doroteo Torre Toledano	UAM, Spain
Eduardo Lleida	UZ, Spain
Eric Laporte	Université Paris Est, France
Eva Navas	UBC, Spain
Fabio Kepler	USP, Brazil
Fábio Violaro	Unicamp, Brazil
Fernando Resende	UFRJ, Brazil
Gaël Harry Dias	UBI, Portugal
Gladis Almeida	UFSCAR, Brazil
Helena de Medeiros Caseli	UFSCAR, Brazil
Irene Rodrigues	UE, Portugal
Isabel Falé	Universidade Aberta, Portugal
Isabel Trancoso	INESC-ID/IST, Portugal
Ivandré Paraboni	USP, Brazil
Jean-Luc Minel	Université de Paris X, France
João Balsa	UL, Portugal
João Luís Rosa	USP-SC, Brazil
João Paulo Neto	INESC-ID/IST, Portugal
João Veloso	UP, Portugal
Joaquim Ferreira da Silva	UNL, Portugal
Joaquim Llisterri	UAB, Spain

Jorge Baptista	Universidade do Algarve, Portugal
José Gabriel Lopes	UNL, Portugal
José João Almeida	UM, Portugal
Julia Hirschberg	Columbia University, USA
Laura Alonso Alemany	University National of Cordoba, Argentina
Leandro Oliveira	Embrapa, Brazil
Leandro Wives	UFRGS, Brazil
Lúcia Rino	UFSCAR, Brazil
Lucia Specia	University of Wolverhampton, UK
Luís Oliveira	INESC-ID, Portugal
Luís Sá	UC, Portugal
Luísa Coheur	INESC-ID/IST, Portugal
Luiz Pizzato	University of Sydney, Australia
Magali Duran	USP-SC, Brazil
Mara Abel	UFRGS, Brazil
Marcelo Finger	USP, Brazil
Marco Gonzalez	PUC-RS, Brazil
Maria das Graças Volpe Nunes	USP-SC, Brazil
Maria Helena Mira Mateus	ILTEC, Portugal
Maria José Finatto	UFRGS, Brazil
Mário Silva	INESC-ID/IST, Portugal
Michel Gagnon	Ecole Polytechnique, Canada
Miguel Sales Dias	Microsoft-MLDC, Portugal
Nuno Cavalheiro Marques	UNL, Portugal
Nuno Mamede	INESC-ID/IST, Portugal
Pablo Gamallo	University of Santiago de Compostela, Spain
Palmira Marrafa	UL, Portugal
Paulo Gomes	UC, Portugal
Paulo Quaresma	UE, Portugal
Plínio Barbosa	Unicamp, Brazil
Ranniery Maia	Toshiba, UK
Renata Vieira	PUC-RS, Brazil
Ricardo Ribeiro	INESC-ID/ISCTE-IUL, Portugal
Robert Dale	Macquarie University, Australia
Ronaldo Martins	Univas, Brazil
Rove Chishman	Unisinos, Brazil
Rubén San-Segundo	UPM, Spain
Ruy Luiz Milidiú	PUC-Rio, Brazil
Sandra Aluisio	USP-SC, Brazil
Sergio Freitas	UnB, Brazil
Solange Rezende	USP-SC, Brazil
Stanley Loh	UCPEL, Brazil
Steven Bird	University of Melbourne, Australia
Thiago Pardo	USP-SC, Brazil
Tracy Holloway King	Microsoft, USA

Valéria Feltrim	UEM, Brazil
Vera Strube de Lima	PUC-RS, Brazil
Violeta Quental	PUC-Rio, Brazil
Vitor Rocio	Universidade Aberta, Portugal
Viviane Moreira	UFRGS, Brazil

Steering Committee

António Teixeira	Universidade de Aveiro, Portugal (Chair)
Fernando Perdigão	Universidade de Coimbra, Portugal
Jorge Baptista	Universidade do Algarve, Portugal
Renata Vieira	PUC-RS, Brazil
Violeta Quental	PUC-RJ, Brazil

Table of Contents

Phonology, Morphology and POS-Tagging

Verb Analysis in a Highly Inflective Language with an MFF Algorithm.....	1
<i>António Branco and Filipe Nunes</i>	
Automatic Analysis of Portuguese Verb Morphology: Solving Ambiguities Caused by Thematic Vowel Allomorphs	12
<i>Vera Vasilévski, Leonor Schiar-Cabral, and Márcio José Araújo</i>	
Coordination of <i>-mente</i> Ending Adverbs in Portuguese: An Integrated Solution	24
<i>Jorge Baptista, Lucas Nunes Vieira, Cláudio Diniz, and Nuno Mamede</i>	
Morphosyntactic Analysis of Language in Children with Autism Spectrum Disorder	35
<i>Raquel Reis and António Teixeira</i>	
<i>Lince</i> , an End User Tool for the Implementation of the Spelling Reform of Portuguese.....	46
<i>José Pedro Ferreira, António Lourinho, and Margarita Correia</i>	
Searching a Mixed Corpus in the Light of the New Portuguese Orthographic Norm	56
<i>Gracinda Carvalho, Isabel Falé, David Martins de Matos, and Vitor Rocio</i>	

Acquisition

Extraction of Bilingual Cognates from Wikipedia	63
<i>Pablo Gamallo and Marcos Garcia</i>	
Corpus-Based Acquisition of Support Verb Constructions for Portuguese	73
<i>Britta D. Zeller and Sebastian Padó</i>	
Improving Portuguese Term Extraction	85
<i>Lucelene Lopes and Renata Vieira</i>	
A Method for Automatically Extracting Domain Semantic Networks from Wikipedia	93
<i>Clarissa Castellã Xavier and Vera Lúcia Strube de Lima</i>	

Extracting Temporal Information from Portuguese Texts	99
<i>Francisco Costa and António Branco</i>	

It Is the Time for Portuguese Texts!	106
<i>Olga Craveiro, Joaquim Macedo, and Henrique Madeira</i>	

Language Resources

A Large Portuguese Corpus On-Line: Cleaning and Preprocessing	113
<i>Michel Génèreux, Iris Hendrickx, and Amália Mendes</i>	

Dicionário-Aberto: A Source of Resources for the Portuguese Language Processing	121
<i>Alberto Simões, Álvaro Iriarte Sanromán, and José João Almeida</i>	

Towards a Common Sense Base in Portuguese for the Linked Open Data Cloud	128
<i>Glória Pinheiro, Vasco Furtado, Tarcisio Pequeno, and Caio Ferreira</i>	

Linguistic Description, Syntax and Parsing

Weak Object Pronouns in Brazilian Portuguese: An LFG Analysis	139
<i>Ana R. Luís</i>	

Entropy-Guided Feature Generation for Structured Learning of Portuguese Dependency Parsing	146
<i>Eraldo R. Fernandes and Ruy L. Milidiú</i>	

Bayesian Induction of Syntactic Language Models for Brazilian Portuguese	157
<i>Daniel Emilio Beck and Helena de Medeiros Caseli</i>	

Automatic Generation of Cloze Question Stems	168
<i>Rui Correia, Jorge Baptista, Maxine Eskenazi, and Nuno Mamede</i>	

Semantics

Toponym Disambiguation Using Ontology-Based Semantic Similarity	179
<i>David S. Batista, João D. Ferreira, Francisco M. Couto, and Mário J. Silva</i>	

Automatic Hyponymy Identification from Brazilian Portuguese Texts	186
<i>Leonardo Sameshima Taba and Helena de Medeiros Caseli</i>	

Semantic Role Labeling for Portuguese – A Preliminary Approach –	193
<i>João Sequeira, Teresa Gonçalves, and Paulo Quaresma</i>	

An Architecture for Semantic Role Labeling on Portuguese	204
<i>Erick Rocha Fonseca and João Luís G. Rosa</i>	

Towards Semi-supervised Brazilian Portuguese Semantic Role Labeling: Building a Benchmark	210
<i>Fernando Emilio Alva-Mancheogo and João Luís G. Rosa</i>	

Opinion Analysis

Building a Sentiment Lexicon for Social Judgement Mining	218
<i>Mário J. Silva, Paula Carvalho, and Luís Sarmento</i>	

A Bootstrapping Algorithm for Learning the Polarity of Words	229
<i>António Paulo Santos, Hugo Gonçalves Oliveira, Carlos Ramos, and Nuno C. Marques</i>	

The Role of Language Registers in Polarity Propagation	235
<i>António Paulo Santos, Hugo Gonçalves Oliveira, Carlos Ramos, and Nuno C. Marques</i>	

Sentiment Analysis on Twitter Data for Portuguese Language	241
<i>Marlo Souza and Renata Vieira</i>	

Natural Language Processing Applications

REAP.PT Serious Games for Learning Portuguese	248
<i>André Silva, Cristiano Marques, Jorge Baptista, Alfredo Ferreira Jr., and Nuno Mamede</i>	

Graph-Based Methods for Multi-document Summarization: Exploring Relationship Maps, Complex Networks and Discourse Information	260
<i>Rafael Ribaldo, Ademar Takeo Akabane, Lucia Helena Machado Rino, and Thiago Alexandre Salgueiro Pardo</i>	

Clustering and Categorization of Brazilian Portuguese Legal Documents	272
<i>Luís Otávio de Colla Furquim and Vera Lúcia Strube de Lima</i>	

SIGA, a System to Manage Information Retrieval Evaluations	284
<i>Luís Costa, Cristina Mota, and Diana Santos</i>	

E-commerce Market Analysis from a Graph-Based Product Classifier . . .	291
<i>Andréa Britto Mattos, Marcelo Van Kampen, Camila Carriço, André Ricardo Dias, and Alexandre Crivellaro</i>	

A Description Logic for InferenceNet.Br	298
<i>Wellington Franco, Thiago Alves, Henrique Viana, and João Alcântara</i>	

Speech Production and Phonetics

Real-Time MRI for Portuguese Database, Methods and Applications ...	306
<i>António Teixeira, Paula Martins, Catarina Oliveira, Carlos Ferreira, Augusto Silva, and Ryan Shosted</i>	
Production and Modeling of the European Portuguese Palatal Lateral	318
<i>António Teixeira, Paula Martins, Catarina Oliveira, and Augusto Silva</i>	
A New Methodology for Comparing Speech Rhythm Structure between Utterances: Beyond Typological Approaches	329
<i>Plínio A. Barbosa and Wellington da Silva</i>	
Constructing Physically Realistic VCV Stimuli for the Perception of Stop Voicing in European Portuguese	338
<i>Daniel Pape, Luis M.T. Jesus, and Pascal Perrier</i>	
Automatic Phonetic Transcription by Phonological Derivation	350
<i>Marcos Garcia and Isaac J. González</i>	

Speech Resources

The C-ORAL-BRASIL I: Reference Corpus for Informal Spoken Brazilian Portuguese	362
<i>Tommaso Raso and Heliana Mello</i>	
A European Portuguese Children Speech Database for Computer Aided Speech Therapy	368
<i>Carla Lopes, Arlindo Veiga, and Fernando Perdigão</i>	
Baseline Acoustic Models for Brazilian Portuguese Using CMU Sphinx Tools	375
<i>Rafael Oliveira, Pedro Batista, Nelson Neto, and Aldebaro Klautau</i>	

Speech Processing and Applications

A Fishervoice-SVM Language Identification System	381
<i>Paula Lopez-Otero, Laura Docio-Fernandez, and Carmen Garcia-Mateo</i>	
Summarizing Speech by Contextual Reinforcement of Important Passages	392
<i>Ricardo Ribeiro and David Martins de Matos</i>	
Incorporating ASR Information in Spoken Dialog System Confidence Score	403
<i>José Lopes, Marine Eskenazi, and Isabel Trancoso</i>	

Transcription of Multi-variety Portuguese Media Contents	409
<i>Alberto Abad, Hugo Meinedo, Isabel Trancoso, and João Neto</i>	
Towards Automatic Classification of Speech Styles	421
<i>Arlindo Veiga, Sara Candeias, Dirce Celorico, Jorge Proença, and Fernando Perdigão</i>	
Author Index	427