## Lecture Notes in Artificial Intelligence 11658

## Subseries of Lecture Notes in Computer Science

### Series Editors

Randy Goebel
University of Alberta, Edmonton, Canada

Yuzuru Tanaka
Hokkaido University, Sapporo, Japan

Wolfgang Wahlster
DFKI and Saarland University, Saarbrücken, Germany

## Founding Editor

Jörg Siekmann

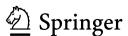
DFKI and Saarland University, Saarbrücken, Germany

More information about this series at http://www.springer.com/series/1244

Albert Ali Salah · Alexey Karpov · Rodmonga Potapova (Eds.)

# Speech and Computer

21st International Conference, SPECOM 2019 Istanbul, Turkey, August 20–25, 2019 Proceedings



Editors
Albert Ali Salah
Utrecht University
Utrecht, The Netherlands

Boğaziçi University Istanbul, Turkey

Rodmonga Potapova 
Moscow State Linguistic University
Moscow, Russia

Alexey Karpov St. Petersburg Institute for Informatics and Automation of the Russian Academy of Sciences
St. Petersburg, Russia

ISSN 0302-9743 ISSN 1611-3349 (electronic) Lecture Notes in Artificial Intelligence ISBN 978-3-030-26060-6 ISBN 978-3-030-26061-3 (eBook) https://doi.org/10.1007/978-3-030-26061-3

LNCS Sublibrary: SL7 - Artificial Intelligence

#### © Springer Nature Switzerland AG 2019

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, 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.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

## SPECOM 2019 Preface

The International Conference on Speech and Computer (SPECOM) was established by the St. Petersburg Institute for Informatics and Automation of the Russian Academy of Sciences (SPIIRAS) and the Herzen State Pedagogical University of Russia thanks to the efforts of Prof. Yuri Kosarev and Prof. Rajmund Piotrowski.

In its long history, the SPECOM conference was organized alternately by SPIIRAS and by the Moscow State Linguistic University (MSLU) in their home cities. SPECOM 2019 was the 21st event in the series, organized by Boğaziçi University (Istanbul, Turkey), in cooperation with SPIIRAS and MSLU. The conference was sponsored by ASM Solutions Ltd. (Moscow, Russia) and supported by the International Speech Communication Association. The conference was held jointly with the 4th International Conference on Interactive Collaborative Robotics (ICR) – where problems and modern solutions of human–robot interaction were discussed – during August 20–25, 2019 at Boğaziçi University, one of the top research universities in Turkey, established in 1863.

During the conferences three invited talks were given by Prof. Hynek Hermansky (Julian S. Smith Professor of Electrical Engineering and the Director of the Center for Language and Speech Processing at the Johns Hopkins University in Baltimore, Maryland, USA and Research Professor at the Brno University of Technology, Czech Republic), Prof. Odette Scharenborg (Delft University of Technology, The Netherlands), and Prof. Erol Şahin (Computer Engineering Dept., Middle East Technical University, Ankara, Turkey).

It is often argued that in in processing of sensory signals such as speech, engineering should apply knowledge of properties of human perception—both have the same goal of getting information from the signal. Prof. Hermansky's talk, entitled "If You Can't Beat Them, Join Them," showed examples from speech technology that perceptual research can also learn from advances in technology. Since speech evolved to be heard and properties of hearing are imprinted on speech, engineering optimizations of speech technology often yield human-like processing strategies. Prof. Hermansky presented a model of human speech communication which suggests that redundancies introduced in speech production in order to protect the message during its transmission through a realistic noisy acoustic environment are being used by human speech perception for a reliable decoding of the message. That led to a particular architecture of an automatic recognition (ASR) system in which longer temporal segments of spectrally smoothed temporal trajectories of spectral energies in individual frequency bands of speech are used to derive estimates of the posterior probabilities of speech sounds. Combinations of these estimates in reliable frequency bands were then adaptively fused to yield the final probability vectors, which best satisfy the adopted performance monitoring criteria.

Speech recognition is the mapping of a continuous, highly variable speech signal onto discrete, abstract representations. In both human and automatic speech processing, the phoneme is considered to play an important role. Abstractionist theories of human

speech processing assume the presence of abstract, phoneme-like units that sequenced together constitute words, while many large vocabulary automatic speech recognition (ASR) systems use phoneme acoustic models. Prof. Scharenborg, in her talk entitled "The Representation of Speech in Human and Artificial Brain," argued that phonemes might not be the unit of speech representation during human speech processing and that comparisons between humans and dynamic neural networks and cross-fertilization of the two research fields can provide valuable insights into the way humans process speech and thereby improve ASR technology. The present volume includes an invited paper by Prof. Scharenborg that discusses these issues at length.

Prof. Şahin's talk, entitled "Animating Industrial Robots for Human–Robot Interaction," discussed interesting interaction research for robot-assisted assembly operations in the production lines of factories. Prof. Şahin argued that as platforms that require fast and fine manipulation of parts and tools remain beyond the capabilities of robotic systems in the near future, robotic systems are predicted not to replace, but to collaborate with the humans working on the assembly lines to increase their productivity. He briefly summarized the vision and goals of a recent TUBITAK project, titled CIRAK, which aims to develop a robotic manipulator system that will help humans in an assembly task by handing them the proper tools and parts at the right time in a correct manner. Toward this end, Prof. Şahin shared his group's recent studies on the creation of a commercial robotic manipulator platform, made more life-like by extensions and modifications to its look and behavior.

This volume contains a collection of submitted papers presented at the conference, which were thoroughly reviewed by members of the Program Committee consisting of more than 100 top specialists, as well as an invited paper by Prof. Scharenborg. Each paper was reviewed, single blind, by two to four committee members (three reviewers on the average) and then discussed by the program chairs. In total, 57 papers were selected by the Program Committee for presentation at the SPECOM Conference. A total of 126 submissions were received and evaluated for SPECOM/ICR. The conference sessions were thematically organized, into Audio Signal Processing, Automatic Speech Recognition, Speaker Recognition, Computational Paralinguistics, Speech Synthesis, Sign Language and Multimodal Processing, and Speech and Language Resources. An increasing number of papers used deep neural network-based approaches across these themes.

We would like to express our gratitude to all authors for providing their papers on time, to the members of the Program Committee for their careful reviews and paper selection, and to the editors and correctors for their hard work in preparing this volume. Special thanks are due to Alen Demirel, Cem Tunçel, Bilge Yüksel, Hasan Küçük of BROS Group, our Conference Office, for their excellent work during the conference organization.

August 2019

Albert Ali Salah Alexey Karpov Rodmonga Potapova

## **Organization**

#### General Co-chairs

Albert Ali Salah Utrecht University and Boğazici University,

The Netherlands/Turkey

SPIIRAS Institute, Russia Alexey Karpov

Rodmonga Potapova Moscow State Linguistic University, Russia

## **Program Co-chairs**

Heysem Kaya Namık Kemal University, Turkey Murat Saraclar Boğaziçi University, Turkey Ebru Arisov Saraclar MEF University, Turkey

## **Program Committee**

Shyam Agrawal, India Ruediger Hoffmann, Germany Lale Akarun, Turkey Marek Hruz, Czech Republic Rafet Akdeniz, Turkey Kazuki Irie, Germany

Tanel Alumäe, Estonia Levent M. Arslan, Turkey Ebru Arisov Saraclar, Turkey

Elias Azarov, Belarus

Peter Beim Graben, Germany Marie-Luce Bourguet, UK

Christian-Alexander Bunge, Germany

Eric Castelli, Vietnam Vladimir Chuchupal, Russia Nicholas Cummins, Germany

Vlado Delic, Serbia Hasan Demir, Turkey Olivier Deroo, Belgium Anna Esposito, Italy Keelan Evanini, USA Vera Evdokimova, Russia Nikos Fakotakis, Greece Mauro Falcone, Italy

Vasiliki Foufi, Switzerland

Philip N. Garner, Switzerland Gábor Gosztolya, Hungary

Tunga Gungor, Turkey Ivan Himawan, Australia

Rainer Jaeckel, Germany Oliver Jokisch, Germany Denis Jouvet, France Alexey Karpov, Russia Heysem Kaya, Turkey Andreas Kerren, Sweden Tomi Kinnunen, Finland Irina Kipyatkova, Russia Daniil Kocharov, Russia Liliya Komalova, Russia Evgeny Kostyuchenko, Russia Ivan Kraljevski, Germany

Galina Lavrentyeva, Russia Benjamin Lecouteux, France Boris Lobanov, Belarus Elena Lyakso, Russia Joseph Mariani, France Maria De Marsico, Italy

Jindřich Matoušek, Czech Republic

Yuri Matveev, Russia

Li Meng, UK

Peter Mihajlik, Hungary Iosif Mporas, UK

#### Organization

viii

Bernd Möbius, Germany Luděk Müller, Czech Republic Satoshi Nakamura, Japan Stavros Ntalampiras, Italy Géza Németh, Hungary Olga Perepelkina, Russia Dimitar Popov, Bulgaria Branislav Popović, Serbia Rodmonga Potapova, Russia Fabien Ringeval, France Andrey Ronzhin, Russia Paolo Rosso, Spain Sakriani Sakti, Japan Albert Ali Salah, The Netherlands/Turkey Murat Saraçlar, Turkey Maximilian Schmitt, Germany Friedhelm Schwenker, Germany Vidhyasaharan Sethu, Australia

Ingo Siegert, Germany Vered Silber-Varod, Israel Pavel Skrelin, Russia Mikhail Stolbov, Russia Tilo Strutz, Germany György Szaszák, Hungary Ivan Tashev, USA Natalia Tomashenko, France Laszlo Toth, Hungary Isabel Trancoso, Portugal Jan Trmal, USA Liliya Tsirulnik, USA Dirk Van Compernolle, Belgium Vasilisa Verkhodanova, The Netherlands Benjamin Weiss, Germany Andreas Wendemuth, Germany Matthias Wolff, Germany Milos Zelezny, Czech Republic Zixing Zhang, UK

### **Additional Reviewers**

Milan Sečujski, Serbia

Cem Rıfkı Aydın Somnath Banerjee Branko Brkljač Buse Buz Koray Çiftçi Gretel Liz De la Peña Sarracén Ali Erkan Nikša Jakovljević Atul Kumar Alexander Leipnitz Michael Maruschke Iris Ouyang Sergey Rybin Siniša Suzić Juan Javier Sánchez Junquera Oxana Verkholyak Celaleddin Yeroglu

## **Contents**

The Representation of Speech and Its Processing in the Human Brain and Deep Neural Networks	1
A Detailed Analysis and Improvement of Feature-Based Named Entity Recognition for Turkish	9
A Comparative Study of Classical and Deep Classifiers for Textual Addressee Detection in Human-Human-Machine Conversations	20
Acoustic Event Mixing to Multichannel AMI Data for Distant Speech Recognition and Acoustic Event Classification Benchmarking  Sergei Astapov, Gleb Svirskiy, Aleksandr Lavrentyev, Tatyana Prisyach, Dmitriy Popov, Dmitriy Ubskiy, and Vladimir Kabarov	31
Speech-Based L2 Call System for English Foreign Speakers	43
A Pattern Mining Approach in Feature Extraction for Emotion Recognition from Speech	54
Towards a Dialect Classification in German Speech Samples	64
Classification of Regional Accent Using Speech Rhythm Metrics	75
PocketEAR: An Assistive Sound Classification System for Hearing-Impaired	82
Time-Continuous Emotion Recognition Using Spectrogram Based CNN-RNN Modelling	93
Developmental Disorders Manifestation in the Characteristics of the Child's Voice and Speech: Perceptual and Acoustic Study	103

Lenar Gabdrakhmanov, Rustem Garaev, and Evgenii Razinkov	113
Differentiating Laughter Types via HMM/DNN and Probabilistic Sampling	122
Gábor Gosztolya, András Beke, and Tilda Neuberger	
Word Discovering in Low-Resources Languages Through Cross-Lingual Phonemes	133
Semantic Segmentation of Historical Documents via Fully-Convolutional Neural Network	142
A New Approach of Adaptive Filtering Updating for Acoustic Echo Cancellation	150
Code-Switching Language Modeling with Bilingual Word Embeddings:  A Case Study for Egyptian Arabic-English	160
Identity Extraction from Clusters of Multi-modal Observations	171
Don't Talk to Noisy Drones – Acoustic Interaction with Unmanned Aerial Vehicles	180
Method for Multimodal Recognition of One-Handed Sign Language Gestures Through 3D Convolution and LSTM Neural Networks  Ildar Kagirov, Dmitry Ryumin, and Alexandr Axyonov	191
LSTM-Based Kazakh Speech Synthesis	201
Combination of Positions and Angles for Hand Pose Estimation	209
LSTM-Based Language Models for Very Large Vocabulary Continuous Russian Speech Recognition System	219

Svarabhakti Vowel Occurrence and Duration in Rhotic Clusters in French Lyric Singing	227
The Evaluation Process Automation of Phrase and Word Intelligibility Using Speech Recognition Systems	237
Detection of Overlapping Speech for the Purposes of Speaker Diarization Marie Kunešová, Marek Hrúz, Zbyněk Zajíc, and Vlasta Radová	247
Exploring Hybrid CTC/Attention End-to-End Speech Recognition with Gaussian Processes.  Ludwig Kürzinger, Tobias Watzel, Lujun Li, Robert Baumgartner, and Gerhard Rigoll	258
Estimating Aggressiveness of Russian Texts by Means of Machine Learning	270
Software Subsystem Analysis of Prosodic Signs of Emotional Intonation Boris Lobanov and Vladimir Zhitko	280
Assessing Alzheimer's Disease from Speech Using the i-vector Approach José Vicente Egas López, László Tóth, Ildikó Hoffmann, János Kálmán, Magdolna Pákáski, and Gábor Gosztolya	289
AD-Child.Ru: Speech Corpus for Russian Children with Atypical Development	299
Building a Pronunciation Dictionary for the Kabyle Language	309
Speech-Based Automatic Assessment of Question Making Skill in L2 Language	317
Automatic Recognition of Speaker Age and Gender Based on Deep  Neural Networks	327
Investigating Joint CTC-Attention Models for End-to-End Russian  Speech Recognition	337

Author Clustering with and Without Topical Features	348
Assessment of Syllable Intelligibility Based on Convolutional Neural Networks for Speech Rehabilitation After Speech Organs Surgical Interventions  Evgeny Kostuchenko, Dariya Novokhrestova, Svetlana Pekarskikh, Alexander Shelupanov, Mikhail Nemirovich-Danchenko, Evgeny Choynzonov, and Lidiya Balatskaya	359
Corpus Study of Early Bulgarian Onomatopoeias in the Terms of CHILDES	370
EEG Investigation of Brain Bioelectrical Activity (Regarding Perception of Multimodal Polycode Internet Discourse)	381
Some Peculiarities of Internet Multimodal Polycode Corpora Annotation Rodmonga Potapova, Vsevolod Potapov, Liliya Komalova, and Andrey Dzhunkovskiy	392
New Perspectives on Canadian English Digital Identity Based on Word Stress Patterns in Lexicon and Spoken Corpus	401
Automatic Speech Recognition for Kreol Morisien: A Case Study for the Health Domain.  Nuzhah Gooda Sahib-Kaudeer, Baby Gobin-Rahimbux, Bibi Saamiyah Bahsu, and Maryam Farheen Aasiyah Maghoo	414
Script Selection Using Convolutional Auto-encoder for TTS  Speech Corpus	423
Pragmatic Markers Distribution in Russian Everyday Speech: Frequency Lists and Other Statistics for Discourse Modeling	433
Curriculum Learning in Sentiment Analysis	444
First Minute Timing in American Telephone Talks:  A Cognitive Approach	451

Contents	xiii
Syntactic Segmentation of Spontaneous Speech: Psychological and Cognitive Aspects	459
Dual-Microphone Speech Enhancement System Attenuating both Coherent and Diffuse Background Noise	471
Reducing the Inter-speaker Variance of CNN Acoustic Models Using Unsupervised Adversarial Multi-task Training	481
Estimates of Transmission Characteristics Related to Perception of Bone-Conducted Speech Using Real Utterances and Transcutaneous Vibration on Larynx	491
Singing Voice Database	501
How Dysarthric Prosody Impacts Naïve Listeners' Recognition Vass Verkhodanova, Sanne Timmermans, Matt Coler, Roel Jonkers, Bauke de Jong, and Wander Lowie	510
Light CNN Architecture Enhancement for Different Types Spoofing Attack Detection	520
Deep Neural Network Quantizers Outperforming Continuous Speech Recognition Systems	530
Speaking Style Based Apparent Personality Recognition	540
Diarization of the Language Consulting Center Telephone Calls Zbyněk Zajíc, Josef V. Psutka, Lucie Zajícová, Luděk Müller, and Petr Salajka	549
NN-Based Czech Sign Language Synthesis	559
Re-evaluation of Words Used in Speech Audiometry	569
Author Index	579