

24th Annual Conference of the International Speech Communication Association (INTERSPEECH 2023)

Dublin, Ireland
20-24 August 2023

Volume 1 of 8

ISBN: 978-1-7138-8880-2

Printed from e-media with permission by:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571



Some format issues inherent in the e-media version may also appear in this print version.

Copyright© (2023) by International Speech Communication Association
All rights reserved.

Printed with permission by Curran Associates, Inc. (2024)

For permission requests, please contact International Speech Communication Association
at the address below.

International Speech Communication Association
c/o Mme Emmanuelle FOXONET
4 Rue des Fauvettes - Lous Tourils
F-66390 Baixas, France

Phone: 49 228 735 643
Fax: 33 468 385 827

secretariat@isca-speech.org

Additional copies of this publication are available from:

Curran Associates, Inc.
57 Morehouse Lane
Red Hook, NY 12571 USA
Phone: 845-758-0400
Fax: 845-758-2633
Email: curran@proceedings.com
Web: www.proceedings.com

TABLE OF CONTENTS

VOLUME 1

KEYNOTE 1 ISCA MEDALLIST

Bridging Speech Science and Technology — Now and into the Future.....	1
<i>Shrikanth Narayanan</i>	

SPEECH SYNTHESIS: PROSODY AND EMOTION

Emotional Talking Head Generation Based on Memory-Sharing and Attention-Augmented Networks	2
<i>Jianrong Wang, Yaxin Zhao, Li Liu, Tianyi Xu, Qi Li, Sen Li</i>	
Speech Synthesis with Self-Supervisedly Learnt Prosodic Representations	7
<i>Zhao-Ci Liu, Zhen-Hua Ling, Ya-Jun Hu, Jia Pan, Jin-Wei Wang, Yun-Di Wu</i>	
EmoMix: Emotion Mixing Via Diffusion Models for Emotional Speech Synthesis	12
<i>Haobin Tang, Xulong Zhang, Jianzong Wang, Ning Cheng, Jing Xiao</i>	
Laughter Synthesis Using Pseudo Phonetic Tokens with a Large-Scale In-The-Wild Laughter Corpus	17
<i>Detai Xin, Shinnosuke Takamichi, Ai Morimatsu, Hiroshi Saruwatari</i>	
Explicit Intensity Control for Accented Text-To-Speech.....	22
<i>Rui Liu, Haolin Zuo, De Hu, Guanglai Gao, Haizhou Li</i>	
Comparing Normalizing Flows and Diffusion Models for Prosody and Acoustic Modelling in Text-To-Speech.....	27
<i>Guangyan Zhang, Thomas Merritt, Sam Ribeiro, Biel Tura-Vecino, Kayoko Yanagisawa, Kamil Pokora, Abdelhamid Ezzerg, Sebastian Cygert, Ammar Abbas, Piotr Bilinski, Roberto Barra-Chicote, Daniel Korzekwa, Jaime Lorenzo-Trueba</i>	

STATISTICAL MACHINE TRANSLATION

Modular Speech-To-Text Translation for Zero-Shot Cross-Modal Transfer	32
<i>Paul-Ambroise Duquenne, Holger Schwenk, Benoît Sagot</i>	
Improving Isochronous Machine Translation with Target Factors and Auxiliary Counters	37
<i>Proyag Pal, Brian Thompson, Yogesh Virkar, Prashant Mathur, Alexandra Chronopoulou, Marcello Federico</i>	
StyleS2ST: Zero-Shot Style Transfer for Direct Speech-To-Speech Translation	42
<i>Kun Song, Yi Ren, Yi Lei, Chunfeng Wang, Kun Wei, Lei Xie, Xiang Yin, Zejun Ma</i>	
Joint Speech Translation and Named Entity Recognition.....	47
<i>Marco Gaido, Sara Papi, Matteo Negri, Marco Turchi</i>	
Analysis of Acoustic Information in End-To-End Spoken Language Translation.....	52
<i>Gerard Sant, Carlos Escolano</i>	

LAMASSU: A Streaming Language-Agnostic Multilingual Speech Recognition and Translation Model Using Neural Transducers	57
<i>Peidong Wang, Eric Sun, Jian Xue, Yu Wu, Long Zhou, Yashesh Gaur, Shujie Liu, Jinyu Li</i>	

SELF-SUPERVISED LEARNING IN ASR

DPHuBERT: Joint Distillation and Pruning of Self-Supervised Speech Models	62
<i>Yifan Peng, Yui Sudo, Shakeel Muhammad, Shinji Watanabe</i>	

Automatic Data Augmentation for Domain Adapted Fine-Tuning of Self-Supervised Speech Representations	67
<i>Salah Zaiem, Titouan Parcollet, Slim Essid</i>	

Dual Acoustic Linguistic Self-Supervised Representation Learning for Cross-Domain Speech Recognition	72
<i>Zhao Yang, Dianwen Ng, Chong Zhang, Xiao Fu, Rui Jiang, Wei Xi, Yukun Ma, Chongjia Ni, Eng Siong Chng, Bin Ma, Jizhong Zhao</i>	

O-1: Self-Training with Oracle and 1-Best Hypothesis.....	77
<i>Murali Karthick Baskar, Andrew Rosenberg, Bhuvana Ramabhadran, Kartik Audhkhasi</i>	

MT4SSL: Boosting Self-Supervised Speech Representation Learning by Integrating Multiple Targets	82
<i>Ziyang Ma, Zhisheng Zheng, Changli Tang, Yujin Wang, Xie Chen</i>	

Comparing Self-Supervised Pre-Training and Semi-Supervised Training for Speech Recognition in Languages with Weak Language Models	87
<i>Léa-Marie Lam-Yee-Mui, Lucas Ondel Yang, Ondrej Klejch</i>	

PROSODY

Chinese EFL Learners' Perception of English Prosodic Focus	92
<i>Xinya Zhang, Ying Chen</i>	

Pitch Accent Variation and the Interpretation of Rising and Falling Intonation in American English	97
<i>Thomas Sostarics, Jennifer Cole</i>	

Tonal Coarticulation as a Cue for Upcoming Prosodic Boundary.....	102
<i>Jianjing Kuang, May Pik Yu Chan, Nari Rhee</i>	

Alignment of Beat Gestures and Prosodic Prominence in German.....	107
<i>Sophie Repp, Lara Muhtz, Johannes Heim</i>	

Creak Prevalence and Prosodic Context in Australian English	112
<i>Hannah White, Joshua Penney, Andy Gibson, Anita Szakay, Felicity Cox</i>	

Speech Reduction: Position Within French Prosodic Structure.....	117
<i>Kübra Bodur, Roxane Bertrand, James S. German, Stéphane Rauzy, Corinne Fredouille, Christine Meunier</i>	

SPEECH PRODUCTION

Transvelar Nasal Coupling Contributing to Speaker Characteristics in Non-Nasal Vowels.....	122
<i>Ziyu Zhu, Yujie Chi, Zhao Zhang, Kiyoshi Honda, Jianguo Wei</i>	

Speech Synthesis from Articulatory Movements Recorded by Real-Time MRI	127
<i>Yuto Otani, Shun Sawada, Hidefumi Ohmura, Kouichi Katsurada</i>	
The ART of Conversation: Measuring Phonetic Convergence and Deliberate Imitation in L2-Speech with a Siamese RNN	132
<i>Zheng Yuan, Aldo Pastore, Dorina De Jong, Hao Xu, Luciano Fadiga, Alessandro D'Ausilio</i>	
Did You See That? Exploring the Role of Vision in the Development of Consonant Feature Contrasts in Children with Cochlear Implants.....	137
<i>James Mahshie, Michael Larsen</i>	

DYSARTHIC SPEECH ASSESSMENT

Automatic Assessments of Dysarthric Speech: The Usability of Acoustic-Phonetic Features	141
<i>Loes Van Bommel, Chiara Pesenti, Xue Wei, Helmer Strik</i>	
Classification of Multi-Class Vowels and Fricatives from Patients Having Amyotrophic Lateral Sclerosis with Varied Levels of Dysarthria Severity	146
<i>Chowdam Venkata Thirumala Kumar, Tanuka Bhattacharjee, Yamini Belur, Atchayaram Nalini, Ravi Yadav, Prasanta Kumar Ghosh</i>	
Parameter-Efficient Dysarthric Speech Recognition Using Adapter Fusion and Householder Transformation	151
<i>Jinzi Qi, Hugo Van Hamme</i>	
Few-Shot Dysarthric Speech Recognition with Text-To-Speech Data Augmentation.....	156
<i>Enno Hermann, Mathew Magimai.-Doss</i>	
Latent Phrase Matching for Dysarthric Speech	161
<i>Dianna Yee, Colin Lea, Jaya Narain, Zifang Huang, Lauren Tooley, Jeffrey P. Bigham, Leah Findlater</i>	
Speech Intelligibility Assessment of Dysarthric Speech by Using Goodness of Pronunciation with Uncertainty Quantification	166
<i>Eun Jung Yeo, Kwanghee Choi, Sunhee Kim, Minhwa Chung</i>	

SPEECH CODING: TRANSMISSION AND ENHANCEMENT

CQNV: A Combination of Coarsely Quantized Bitstream and Neural Vocoder for Low Rate Speech Coding	171
<i>Youqiang Zheng, Li Xiao, Weiping Tu, Yuhong Yang, Ximmeng Xu</i>	
Target Speech Extraction with Conditional Diffusion Model.....	176
<i>Naoyuki Kamo, Marc Delcroix, Tomohiro Nakatani</i>	
Towards Fully Quantized Neural Networks for Speech Enhancement.....	181
<i>Elad Cohen, Hai Victor Habi, Arnon Netzer</i>	
Complex Image Generation SwinTransformer Network for Audio Denoising	186
<i>Youshan Zhang, Jialu Li</i>	

**SPEECH RECOGNITION: SIGNAL PROCESSING, ACOUSTIC MODELING,
ROBUSTNESS, ADAPTATION 1**

Using Text Injection to Improve Recognition of Personal Identifiers in Speech.....	191
<i>Yochai Blau, Rohan Agrawal, Lior Madmony, Gary Wang, Andrew Rosenberg, Zhehuai Chen, Zorik Gekhman, Genady Beryozkin, Parisa Haghani, Bhuvana Ramabhadran</i>	
Investigating Wav2vec2 Context Representations and the Effects of Fine-Tuning, a Case-Study of a Finnish Model.....	196
<i>Tamas Grosz, Yaroslav Getman, Ragheb Al-Ghezi, Aku Rouhe, Mikko Kurimo</i>	
Transformer-Based Speech Recognition Models for Oral History Archives in English, German, and Czech.....	201
<i>Jan Lehecka, Jan Švec, Josef V. Psutka, Pavel Ircing</i>	
Iteratively Improving Speech Recognition and Voice Conversion	206
<i>Mayank Kumar Singh, Naoya Takahashi, Naoyuki Onoe</i>	
LABERT: A Combination of Local Aggregation and Self-Supervised Speech Representation Learning for Detecting Informative Hidden Units in Low-Resource ASR Systems	211
<i>Kavan Fatehi, Ayse Kucukyilmaz</i>	
TranUSR: Phoneme-To-Word Transcoder Based Unified Speech Representation Learning for Cross-Lingual Speech Recognition	216
<i>Hongfei Xue, Qijie Shao, Peikun Chen, Pengcheng Guo, Lei Xie, Jie Liu</i>	
Dual-Mode NAM: Effective Top-K Context Injection for End-To-End ASR	221
<i>Zelin Wu, Tsendsuren Munkhdalai, Pat Rondon, Golan Pundak, Khe Chai Sim, Christopher Li</i>	
GhostRNN: Reducing State Redundancy in RNN with Cheap Operations.....	226
<i>Hang Zhou, Xiaoxu Zheng, Yunhe Wang, Michael Bi Mi, Deyi Xiong, Kai Han</i>	
Task-Agnostic Structured Pruning of Speech Representation Models	231
<i>Haoyu Wang, Siyuan Wang, Wei-Qiang Zhang, Suo Hongbin, Yulong Wan</i>	
Factual Consistency Oriented Speech Recognition	236
<i>Naoyuki Kanda, Takuya Yoshioka, Yang Liu</i>	
Multi-Head State Space Model for Speech Recognition	241
<i>Yassir Fathullah, Chunyang Wu, Yuan Shangguan, Junteng Jia, Wenhan Xiong, Jay Mahadeokar, Chunxi Liu, Yangyang Shi, Ozlem Kalinli, Mike Seltzer, Mark J. F. Gales</i>	
Cascaded Multi-Task Adaptive Learning Based on Neural Architecture Search.....	246
<i>Yingying Gao, Shilei Zhang, Zihao Cui, Chao Deng, Junlan Feng</i>	
Probing Self-Supervised Speech Models for Phonetic and Phonemic Information: A Case Study in Aspiration.....	251
<i>Kinan Martin, Jon Gauthier, Canaan Breiss, Roger Levy</i>	
Selective Biasing with Trie-Based Contextual Adapters for Personalised Speech Recognition Using Neural Transducers.....	256
<i>Philip Harding, Sibong Tong, Simon Wiesler</i>	

ANALYSIS OF SPEECH AND AUDIO SIGNALS 1

Robust Prototype Learning for Anomalous Sound Detection.....	261
<i>Xiao-Min Zeng, Yan Song, Ian McLoughlin, Lin Liu, Li-Rong Dai</i>	
A Multimodal Prototypical Approach for Unsupervised Sound Classification	266
<i>Saksham Singh Kushwaha, Magdalena Fuentes</i>	
Robust Audio Anti-Spoofing with Fusion-Reconstruction Learning on Multi-Order Spectrograms.....	271
<i>Penghui Wen, Kun Hu, Wenxi Yue, Sen Zhang, Wanlei Zhou, Zhiyong Wang</i>	
Adapting Language-Audio Models as Few-Shot Audio Learners	276
<i>Jinhua Liang, Xubo Liu, Haohe Liu, Huy Phan, Emmanouil Benetos, Mark D. Plumbley, Wenwu Wang</i>	
TFECN: Time-Frequency Enhanced ConvNet for Audio Classification.....	281
<i>Mengwei Wang, Zhe Yang</i>	
Resolution Consistency Training on Time-Frequency Domain for Semi-Supervised Sound Event Detection	286
<i>Won-Gook Choi, Joon-Hyuk Chang</i>	
Fine-Tuning Audio Spectrogram Transformer with Task-Aware Adapters for Sound Event Detection	291
<i>Kang Li, Yan Song, Ian McLoughlin, Lin Liu, Jin Li, Li-Rong Dai</i>	
Small Footprint Multi-Channel Network for Keyword Spotting with Centroid Based Awareness	296
<i>Dianwen Ng, Yang Xiao, Jia Qi Yip, Zhao Yang, Biao Tian, Qiang Fu, Eng Siong Chng, Bin Ma</i>	
Few-Shot Class-Incremental Audio Classification Using Adaptively-Refined Prototypes	301
<i>Wei Xie, Yanxiong Li, Qianhua He, Wenchang Cao, Tuomas Virtanen</i>	
Interpretable Latent Space Using Space-Filling Curves for Phonetic Analysis in Voice Conversion	306
<i>Mohammad Hassan Vali, Tom Bäckström</i>	
Topological Data Analysis for Speech Processing	311
<i>Eduard Tulchinskii, Kristian Kuznetsov, Laida Kushnareva, Daniil Cherniavskii, Serguei Barannikov, Irina Piontkovskaya, Sergey Nikolenko, Evgeny Burnaev</i>	
Recycle-And-Distill: Universal Compression Strategy for Transformer-Based Speech SSL Models with Attention Map Reusing and Masking Distillation	316
<i>Kangwook Jang, Sungnyun Kim, Se-Young Yun, Hoirin Kim</i>	
Personalized Acoustic Scene Classification in Ultra-Low Power Embedded Devices Using Privacy-Preserving Data Augmentation.....	321
<i>Timm Koppelman, Semih Agcaer, Rainer Martin</i>	
Background Domain Switch: A Novel Data Augmentation Technique for Robust Sound Event Detection	326
<i>Wei-Cheng Lin, Luca Bondi, Shabnam Ghaffarzadegan</i>	
Joint Prediction of Audio Event and Annoyance Rating in an Urban Soundscape by Hierarchical Graph Representation Learning.....	331
<i>Yuanbo Hou, Siyang Song, Cheng Luo, Andrew Mitchell, Qiaoqiao Ren, Weicheng Xie, Jian Kang, Wenwu Wang, Dick Botteldooren</i>	

Anomalous Sound Detection Using Self-Attention-Based Frequency Pattern Analysis of Machine Sounds	336
<i>Hejing Zhang, Jian Guan, Qiaoxi Zhu, Feiyang Xiao, Youde Liu</i>	
Improving Audio-Text Retrieval Via Hierarchical Cross-Modal Interaction and Auxiliary Captions.....	341
<i>Yifei Xin, Yuexian Zou</i>	
Differential Privacy Enabled Dementia Classification: An Exploration of the Privacy-Accuracy Trade-Off in Speech Signal Data.....	346
<i>Suhas BN, Sarah Rajtmajer, Saeed Abdullah</i>	
Learning Emotional Representations from Imbalanced Speech Data for Speech Emotion Recognition and Emotional Text-To-Speech	351
<i>Shijun Wang, Jón Guðnason, Damian Borth</i>	
Towards Multi-Lingual Audio Question Answering.....	356
<i>Swarup Ranjan Behera, Pailla Balakrishna Reddy, Achyut Mani Tripathi, Megavath Bharadwaj Rathod, Tejesh Karavadi</i>	

SPEECH RECOGNITION: ARCHITECTURE, SEARCH, AND LINGUISTIC COMPONENTS

1

Diacritic Recognition Performance in Arabic ASR	361
<i>Hanan Aldarmaki, Ahmad Ghannam</i>	
Personalization for BERT-Based Discriminative Speech Recognition Rescoring.....	366
<i>Jari Kolehmainen, Yile Gu, Aditya Gourav, Prashanth Gurunath Shivakumar, Ankur Gandhe, Ariya Rastrow, Ivan Bulyko</i>	
On the N-Gram Approximation of Pre-Trained Language Models	371
<i>Aravind Krishnan, Jesujoba O. Alabi, Dietrich Klakow</i>	
Record Deduplication for Entity Distribution Modeling in ASR Transcripts.....	376
<i>Tianyu Huang, Chung Hoon Hong, Carl Wivagg, Kanna Shimizu</i>	
Learning When to Trust Which Teacher for Weakly Supervised ASR	381
<i>Aakriti Agrawal, Milind Rao, Anit Kumar Sahu, Gopinath Chennupati, Andreas Stolcke</i>	
Text-Only Domain Adaptation Using Unified Speech-Text Representation in Transducer	386
<i>Lu Huang, Boyu Li, Jun Zhang, Lu Lu, Zejun Ma</i>	

SPEECH RECOGNITION: TECHNOLOGIES AND SYSTEMS FOR NEW APPLICATIONS 1

Syllable Discovery and Cross-Lingual Generalization in a Visually Grounded, Self-Supervised Speech Model.....	391
<i>Puyuan Peng, Shang-Wen Li, Okko Räsänen, Abdelrahman Mohamed, David Harwath</i>	
Prompting the Hidden Talent of Web-Scale Speech Models for Zero-Shot Task Generalization.....	396
<i>Puyuan Peng, Brian Yan, Shinji Watanabe, David Harwath</i>	
Progress and Prospects for Spoken Language Technology: Results from Five Sexennial Surveys	401
<i>Roger K. Moore, Ricard Marxer</i>	

Acoustic Word Embeddings for Untranscribed Target Languages with Continued Pretraining and Learned Pooling	406
<i>Ramon Sanabria, Ondrej Klejch, Hao Tang, Sharon Goldwater</i>	
CASA-ASR: Context-Aware Speaker-Attributed ASR.....	411
<i>Mohan Shi, Zhihao Du, Qian Chen, Fan Yu, Yangze Li, Shiliang Zhang, Jie Zhang, Li-Rong Dai</i>	
Unsupervised Learning of Discrete Latent Representations with Data-Adaptive Dimensionality from Continuous Speech Streams.....	416
<i>Shun Takahashi, Sakriani Sakti</i>	
AD-TUNING: An Adaptive CHILD-TUNING Approach to Efficient Hyperparameter Optimization of Child Networks for Speech Processing Tasks in the SUPERB Benchmark.....	421
<i>Gaobin Yang, Jun Du, Maokui He, Shutong Niu, Baoxiang Li, Jiakui Li, Chin-Hui Lee</i>	
Distilling Knowledge from Gaussian Process Teacher to Neural Network Student.....	426
<i>Jeremy H. M. Wong, Huayun Zhang, Nancy F. Chen</i>	
Segmental SpeechCLIP: Utilizing Pretrained Image-Text Models for Audio-Visual Learning	431
<i>Saurabhchand Bhati, Jesús Villalba, Laureano Moro-Velazquez, Thomas Thebaud, Najim Dehak</i>	
Towards Hate Speech Detection in Low-Resource Languages: Comparing ASR to Acoustic Word Embeddings on Wolof and Swahili	436
<i>Christiaan Jacobs, Nathanaël Carraz Rakotonirina, Everlyn Asiko Chimoto, Bruce A. Bassett, Herman Kamper</i>	
Mitigating Catastrophic Forgetting for Few-Shot Spoken Word Classification Through Meta-Learning	441
<i>Ruan Van Der Merwe, Herman Kamper</i>	
Online Punctuation Restoration Using ELECTRA Model for Streaming ASR Systems.....	446
<i>Martin Poláček, Petr Cerva, Jindrich Ždánský, Lenka Weingartová</i>	
Language Agnostic Data-Driven Inverse Text Normalization.....	451
<i>Szu-Jui Chen, Debjyoti Paul, Yutong Pang, Peng Su, Xuedong Zhang</i>	
How to Estimate Model Transferability of Pre-Trained Speech Models?	456
<i>Zih-Ching Chen, Chao-Han Huck Yang, Bo Li, Yu Zhang, Nanxin Chen, Shuo-Yiin Chang, Rohit Prabhavalkar, Hung-Yi Lee, Tara Sainath</i>	
Transcribing Speech as Spoken and Written Dual Text Using an Autoregressive Model	461
<i>Mana Ihori, Hiroshi Sato, Tomohiro Tanaka, Ryo Masumura, Saki Mizuno, Nobukatsu Hojo</i>	

LEXICAL AND LANGUAGE MODELING FOR ASR

NoRefER: A Referenceless Quality Metric for Automatic Speech Recognition Via Semi-Supervised Language Model Fine-Tuning with Contrastive Learning	466
<i>Kamer Ali Yuksel, Thiago Castro Ferreira, Golará Javadi, Mohamed Al-Badrashiny, Ahmet Gunduz</i>	
Scaling Laws for Discriminative Speech Recognition Rescoring Models	471
<i>Yile Gu, Prashanth Gurunath Shivakumar, Jari Kolehmainen, Ankur Gandhe, Ariya Rastrow, Ivan Bulyko</i>	

Exploring Energy-Based Language Models with Different Architectures and Training Methods for Speech Recognition	476
<i>Hong Liu, Zhaobiao Lv, Zhijian Ou, Wenbo Zhao, Qing Xiao</i>	
Memory Augmented Lookup Dictionary Based Language Modeling for Automatic Speech Recognition	481
<i>Yukun Feng, Ming Tu, Rui Xia, Chuanzeng Huang, Yuxuan Wang</i>	
Memory Network-Based End-To-End Neural ES-KMeans for Improved Word Segmentation.....	486
<i>Yu Iwamoto, Takahiro Shinozaki</i>	
Retraining-Free Customized ASR for Enharmonic Words Based on a Named-Entity-Aware Model and Phoneme Similarity Estimation	491
<i>Yui Sudo, Kazuya Hata, Kazuhiro Nakadai</i>	

LANGUAGE IDENTIFICATION AND DIARIZATION

Lightweight and Efficient Spoken Language Identification of Long-Form Audio.....	496
<i>Winstead Zhu, Md Iftekhar Tanveer, Yang Janet Liu, Seye Ojumu, Rosie Jones</i>	
End to End Spoken Language Diarization with Wav2vec Embeddings	501
<i>Jagabandhu Mishra, Jayadev N Patil, Amartya Chowdhury, Mahadeva Prasanna</i>	
Efficient Spoken Language Recognition Via Multilabel Classification	506
<i>Oriol Nieto, Zeyu Jin, Franck Derroncourt, Justin Salamon</i>	
Description and Analysis of ABC Submission to NIST LRE 2022.....	511
<i>Pavel Matejka, Anna Silnova, Josef Slavicek, Ladislav Mosner, Oldrich Plhot, Michal Klco, Junyi Peng, Themos Stafylakis, Lukáš Burget</i>	
Exploring the Impact of Pretrained Models and Web-Scraped Data for the 2022 NIST Language Recognition Evaluation	516
<i>Tanel Alumäe, Kunnar Kukk, Viet-Bac Le, Claude Barras, Abdel Messaoudi, Waad Ben Kheder</i>	
Advances in Language Recognition in Low Resource African Languages: The JHU-MIT Submission for NIST LRE22	521
<i>Jesús Villalba, Jonas Borgstrom, Maliha Jahan, Saurabh Kataria, Leibny Paola Garcia, Pedro Torres-Carrasquillo, Najim Dehak</i>	

SPEECH QUALITY ASSESSMENT

DeePMOS: Deep Posterior Mean-Opinion-Score of Speech	526
<i>Xinyu Liang, Fredrik Cumlin, Christian Schüldt, Saikat Chatterjee</i>	
The Role of Formant and Excitation Source Features in Perceived Naturalness of Low Resource Tribal Language TTS: An Empirical Study	531
<i>Ashwini Dasare, Pradyoth Hegde, Supriya Shetty, Deepak K T</i>	
A No-Reference Speech Quality Assessment Method Based on Neural Network with Densely Connected Convolutional Architecture.....	536
<i>Wuxuan Gong, Jing Wang, Yitong Liu, Hongwen Yang</i>	
Probing Speech Quality Information in ASR Systems	541
<i>Bao Thang Ta, Minh Tu Le, Nhat Minh Le, Van Hai Do</i>	

Preference-Based Training Framework for Automatic Speech Quality Assessment Using Deep Neural Network	546
<i>Cheng-Hung Hu, Yusuke Yasuda, Tomoki Toda</i>	

Crowdsourced Data Validation for ASR Training	551
<i>Wannaphong Phatthiyaphaibun, Chompakorn Chaksangchaichot, Thanawin Rakthammanon, Ekapol Chuangsuwanich, Sarana Nutanong</i>	

FEATURE MODELING FOR ASR

Re-Investigating the Efficient Transfer Learning of Speech Foundation Model Using Feature Fusion Methods	556
<i>Zhouyuan Huo, Khe Chai Sim, Dongseong Hwang, Tsendsuren Munkhdalai, Tara Sainath, Pedro M. Mengibar</i>	

Robust Automatic Speech Recognition Via WavAugment Guided Phoneme Adversarial Training.....	561
<i>Gege Qi, Yuefeng Chen, Xiaofeng Mao, Xiaojun Jia, Ranjie Duan, Rong Zhang, Hui Xue</i>	

InterFormer: Interactive Local and Global Features Fusion for Automatic Speech Recognition.....	566
<i>Zhi-Hao Lai, Tian-Hao Zhang, Qi Liu, Xinyuan Qian, Li-Fang Wei, Feng Chen, Song-Lu Chen, Xu-Cheng Yin</i>	

Transductive Feature Space Regularization for Few-Shot Bioacoustic Event Detection	571
<i>Yizhou Tan, Haojun Ai, Shengchen Li, Feng Zhang</i>	

Incorporating L2 Phonemes Using Articulatory Features for Robust Speech Recognition.....	576
<i>Jisung Wang, Haram Lee, Myungwoo Oh</i>	

On the (In)Efficiency of Acoustic Feature Extractors for Self-Supervised Speech Representation Learning	581
<i>Titouan Parcollet, Shucong Zhang, Rogier Van Dalen, Alberto Gil C. P. Ramos, Sourav Bhattacharya</i>	

INTERFACING SPEECH TECHNOLOGY AND PHONETICS

Phonemic Competition in End-To-End ASR Models	586
<i>Louis Ten Bosch, Martijn Bentum, Lou Boves</i>	

Automatic Speaker Recognition with Variation Across Vocal Conditions: A Controlled Experiment with Implications for Forensics	591
<i>Vincent Hughes, Jessica Wormald, Paul Foulkes, Philip Harrison, Finnian Kelly, David Van Der Vloed, Poppy Welch, Chenzi Xu</i>	

Exploring Graph Theory Methods for the Analysis of Pronunciation Variation in Spontaneous Speech	596
<i>Bernhard C. Geiger, Barbara Schuppler</i>	

Automatic Speaker Recognition Performance with Matched and Mismatched Female Bilingual Speech Data.....	601
<i>Bryony Nuttall, Philip Harrison, Vincent Hughes</i>	

SPEECH SYNTHESIS: MULTILINGUALITY

FACTSpeech: Speaking a Foreign Language Pronunciation Using Only Your Native Characters	606
<i>Hong-Sun Yang, Ji-Hoon Kim, Yoon-Cheol Ju, Il-Hwan Kim, Byeong-Yeol Kim, Shuk-Jae Choi, Hyung-Yong Kim</i>	
Cross-Lingual Transfer Learning for Phrase Break Prediction with Multilingual Language Model.....	611
<i>Hoyeon Lee, Hyun-Wook Yoon, Jong-Hwan Kim, Jae-Min Kim</i>	
DSE-TTS: Dual Speaker Embedding for Cross-Lingual Text-To-Speech.....	616
<i>Sen Liu, Yiwei Guo, Chenpeng Du, Xie Chen, Kai Yu</i>	
Generating Multilingual Gender-Ambiguous Text-To-Speech Voices	621
<i>Konstantinos Markopoulos, Georgia Maniati, Georgios Vamvoukakis, Nikolaos Ellinas, Georgios Vardaxoglou, Panos Kakoulidis, Junkwang Oh, Gunu Jho, Inchul Hwang, Aimilios Chalamandaris, Pirros Tsiakoulis, Spyros Raptis</i>	
RAD-MMM: Multilingual Multiaccented Multispeaker Text to Speech.....	626
<i>Rohan Badlani, Rafael Valle, Kevin J. Shih, João Felipe Santos, Siddharth Gururani, Bryan Catanzaro</i>	
Multilingual Context-Based Pronunciation Learning for Text-To-Speech	631
<i>Giulia Comini, Sam Ribeiro, Fan Yang, Heereen Shim, Jaime Lorenzo-Trueba</i>	

SPEECH EMOTION RECOGNITION 1

Personalized Adaptation with Pre-Trained Speech Encoders for Continuous Emotion Recognition	636
<i>Minh Tran, Yufeng Yin, Mohammad Soleymani</i>	
The Importance of Calibration: Rethinking Confidence and Performance of Speech Multi-Label Emotion Classifiers	641
<i>Huang-Cheng Chou, Lucas Goncalves, Seong-Gyun Leem, Chi-Chun Lee, Carlos Busso</i>	
A Preliminary Study on Augmenting Speech Emotion Recognition Using a Diffusion Model	646
<i>Mohammad Ibrahim Malik, Siddique Latif, Raja Jurdak, Björn W. Schuller</i>	
Privacy Risks in Speech Emotion Recognition: A Systematic Study on Gender Inference Attack	651
<i>Basmah Alsenani, Tanaya Guha, Alessandro Vinciarelli</i>	
Episodic Memory for Domain-Adaptable, Robust Speech Emotion Recognition.....	656
<i>James Tavernor, Matthew Perez, Emily Mower Provost</i>	
Stable Speech Emotion Recognition with Head-K-Pooling Loss.....	661
<i>Chaoyue Ding, Jiakui Li, Daoming Zong, Baoxiang Li, Tian-Hao Zhang, Qunyan Zhou</i>	

SHOW AND TELL: HEALTH APPLICATIONS AND EMOTION RECOGNITION

A Personalised Speech Communication Application for Dysarthric Speakers	666
<i>Matthew Gibson, Ievgen Karaulov, Oleksii Zhelo, Filip Jurcicek</i>	
Video Multimodal Emotion Recognition System for Real World Applications	668
<i>Sun-Kyung Lee, Jong-Hwan Kim</i>	

Promoting Mental Self-Disclosure in a Spoken Dialogue System	670
<i>Mahdin Rohmatillah, Bobbi Aditya, Li-Jen Yang, Bryan Gautama Ngo, Willianto Sulaiman, Jen-Tzung Chien</i>	
"Select Language, Modality Or Put on a Mask!" Experiments with Multimodal Emotion Recognition	672
<i>Pawel Bujnowski, Bartlomiej Kuzma, Bartlomiej Paziewski, Jacek Rutkowski, Joanna Marhula, Zuzanna Bordzicka, Piotr Andruszkiewicz</i>	
My Vowels Matter: Formant Automation Tools for Diverse Child Speech	674
<i>Hannah Valentine, Joel Macauslan, Maria Grigos, Marisha Speights</i>	
NEMA: An Ecologically Valid Tool for Assessing Hearing Devices, Advanced Algorithms, and Communication in Diverse Listening Environments.....	676
<i>Nicky Chong-White, Arun Sebastian, Jorge Mejia</i>	
When Words Speak Just as Loudly as Actions: Virtual Agent Based Remote Health Assessment Integrating What Patients Say with What They Do.....	678
<i>Vikram Ramanarayanan, David Pautler, Lakshmi Arbatti, Abhishek Hosamath, Michael Neumann, Hardik Kothare, Oliver Roesler, Jackson Liscombe, Andrew Cornish, Doug Habberstad, Vanessa Richter, David Fox, David Suendermann-Oeft, Ira Shoulson</i>	
Stuttering Detection Application	680
<i>Kowshik Siva Sai Motepalli, Vamshiraghusimha Narasinga, Harsha Pathuri, Hina Khan, Sangeetha Mahesh, Ajish K. Abraham, Anil Kumar Vuppala</i>	
Providing Interpretable Insights for Neurological Speech and Cognitive Disorders from Interactive Serious Games.....	682
<i>Mario Zusag, Laurin Wagner</i>	
Automated Neural Nursing Assistant (ANNA): An Over-The-Phone System for Cognitive Monitoring.....	684
<i>Jacob Solinsky, Raymond Finzel, Martin Michalowski, Serguei Pakhomov</i>	
5G-IoT Cloud Based Demonstration of Real-Time Audio-Visual Speech Enhancement for Multimodal Hearing-Aids	686
<i>Ankit Gupta, Abhijeet Bishnu, Mandar Gogate, Kia Dashtipour, Tughrul Arslan, Ahsan Adeel, Amir Hussain, Tharmalingam Ratnarajah, Mathini Sellathurai</i>	
Towards Two-Point Neuron-Inspired Energy-Efficient Multimodal Open Master Hearing Aid.....	688
<i>Mohsin Raza, Adewale Adetomi, Khubaib Ahmed, Amir Hussain, Tughrul Arslan, Ahsan Adeel</i>	

SPOKEN DIALOG SYSTEMS AND CONVERSATIONAL ANALYSIS 1

FC-MTLF: A Fine- And Coarse-Grained Multi-Task Learning Framework for Cross-Lingual Spoken Language Understanding	690
<i>Xuxin Cheng, Wanshi Xu, Ziyu Yao, Zhihong Zhu, Yaowei Li, Hongxiang Li, Yuexian Zou</i>	

VOLUME 2

C ² A-SLU: Cross and Contrastive Attention for Improving ASR Robustness in Spoken Language Understanding	695
<i>Xuxin Cheng, Ziyu Yao, Zhihong Zhu, Yaowei Li, Hongxiang Li, Yuexian Zou</i>	

Tri-Level Joint Natural Language Understanding for Multi-Turn Conversational Datasets	700
<i>Henry Weld, Sijia Hu, Siqu Long, Josiah Poon, Soyeon Han</i>	
Semantic Enrichment Towards Efficient Speech Representations	705
<i>Gaëlle Laperrière, Ha Nguyen, Sahar Ghannay, Bassam Jabaian, Yannick Estève</i>	
Tensor Decomposition for Minimization of E2E SLU Model Toward On-Device Processing	710
<i>Yosuke Kashiwagi, Siddhant Arora, Hayato Futami, Jessica Huynh, Shih-Lun Wu, Yifan Peng, Brian Yan, Emiru Tsunoo, Shinji Watanabe</i>	
DiffSLU: Knowledge Distillation Based Diffusion Model for Cross-Lingual Spoken Language Understanding	715
<i>Tianjun Mao, Chenghong Zhang</i>	
Integrating Pretrained ASR and LM to Perform Sequence Generation for Spoken Language Understanding	720
<i>Siddhant Arora, Hayato Futami, Yosuke Kashiwagi, Emiru Tsunoo, Brian Yan, Shinji Watanabe</i>	
Contrastive Learning Based ASR Robust Knowledge Selection for Spoken Dialogue System	725
<i>Zhiyuan Zhu, Yusheng Liao, Yu Wang, Yunfeng Guan</i>	
Unsupervised Dialogue Topic Segmentation in Hyperdimensional Space	730
<i>Seongmin Park, Jinkyu Seo, Jihwa Lee</i>	
An Investigation of the Combination of Rehearsal and Knowledge Distillation in Continual Learning for Spoken Language Understanding	735
<i>Umberto Cappellazzo, Daniele Falavigna, Alessio Brutti</i>	
Enhancing New Intent Discovery Via Robust Neighbor-Based Contrastive Learning	740
<i>Zhenhe Wu, Xiaoguang Yu, Meng Chen, Liangqing Wu, Jiahao Ji, Zhoujun Li</i>	
Personalized Predictive ASR for Latency Reduction in Voice Assistants	745
<i>Andreas Schwarz, Di He, Maarten Van Segbroeck, Mohammed Hethnawi, Ariya Rastrow</i>	
Compositional Generalization in Spoken Language Understanding	750
<i>Avik Ray, Yilin Shen, Hongxia Jin</i>	
Sampling Bias in NLU Models: Impact and Mitigation	755
<i>Zefei Li, Anil Ramakrishna, Anna Rumshisky, Andy Rosenbaum, Saleh Solta, Rahul Gupta</i>	
SIDER: Unified Query Rewriting for Steering, Intent Carryover, Disfluencies, Entity Carryover and Repair	760
<i>Jiarui Lu, Bo-Hsiang Tseng, Joel Ruben Antony Moniz, Site Li, Xueyun Zhu, Hong Yu, Murat Akbacak</i>	
Emotion Awareness in Multi-Utterance Turn for Improving Emotion Prediction in Multi-Speaker Conversation	765
<i>Xiaohan Shi, Xingfeng Li, Tomoki Toda</i>	
WhiSLU: End-To-End Spoken Language Understanding with Whisper	770
<i>Minghan Wang, Yinglu Li, Jiaxin Guo, Xiaosong Qiao, Zongyao Li, Hengchao Shang, Daimeng Wei, Shimin Tao, Min Zhang, Hao Yang</i>	

SPEECH CODING AND ENHANCEMENT 1

Biophysically-Inspired Single-Channel Speech Enhancement in the Time Domain.....	775
<i>Chuan Wen, Sarah Verhulst</i>	
On-Device Speaker Anonymization of Acoustic Embeddings for ASR Based on Flexible Location Gradient Reversal Layer.....	780
<i>Md Asif Jalal, Pablo Peso Parada, Jisi Zhang, Mete Ozay, Karthikeyan Saravanan, Myoungji Han, Jung In Lee, Seokyeong Jung</i>	
How to Construct Perfect and Worse-Than-Coin-Flip Spoofing Countermeasures: A Word of Warning on Shortcut Learning.....	785
<i>Hye-Jin Shim, Rosa Gonzalez Hautamäki, Md Sahidullah, Tomi Kinnunen</i>	
CleanUNet 2: A Hybrid Speech Denoising Model on Waveform and Spectrogram.....	790
<i>Zhifeng Kong, Wei Ping, Amrith Dantrey, Bryan Catanzaro</i>	
A Two-Stage Progressive Neural Network for Acoustic Echo Cancellation	795
<i>Zhuangqi Chen, Xianjun Xia, Cheng Chen, Xianke Wang, Yanhong Leng, Li Chen, Roberto Togneri, Yijian Xiao, Piao Ding, Shenyi Song, Pingjian Zhang</i>	
An Intra-BRNN and GB-RVQ Based END-TO-END Neural Audio Codec	800
<i>Linpings Xu, Jiawei Jiang, Dejun Zhang, Xianjun Xia, Li Chen, Yijian Xiao, Piao Ding, Shenyi Song, Sixing Yin, Ferdous Sohel</i>	
Real-Time Personalised Speech Enhancement Transformers with Dynamic Cross-Attended Speaker Representations	804
<i>Shucong Zhang, Malcolm Chadwick, Alberto Gil C. P. Ramos, Titouan Parcollet, Rogier Van Dalen, Sourav Bhattacharya</i>	
CFTNet: Complex-Valued Frequency Transformation Network for Speech Enhancement	809
<i>Nursadul Mamun, John H. L. Hansen</i>	
Feature Normalization for Fine-Tuning Self-Supervised Models in Speech Enhancement	814
<i>Hejung Yang, Hong-Goo Kang</i>	
Multi-Mode Neural Speech Coding Based on Deep Generative Networks.....	819
<i>Wei Xiao, Wenzhe Liu, Meng Wang, Shan Yang, Yupeng Shi, Yuyong Kang, Dan Su, Shidong Shang, Dong Yu</i>	
Streaming Dual-Path Transformer for Speech Enhancement	824
<i>Soo Hyun Bae, Seok Wan Chae, Youngseok Kim, Keunsang Lee, Hyunjin Lim, Lae-Hoon Kim</i>	
Sequence-To-Sequence Multi-Modal Speech In-Painting.....	829
<i>Mahsa Kadkhodaei Elyaderani, Shahram Shirani</i>	
Hybrid AHS: A Hybrid of Kalman Filter and Deep Learning for Acoustic Howling Suppression	834
<i>Hao Zhang, Meng Yu, Yuzhong Wu, Tao Yu, Dong Yu</i>	
Differentially Private Adapters for Parameter Efficient Acoustic Modeling	839
<i>Chun-Wei Ho, Chao-Han Huck Yang, Sabato Marco Siniscalchi</i>	
Incorporating Ultrasound Tongue Images for Audio-Visual Speech Enhancement Through Knowledge Distillation.....	844
<i>Rui-Chen Zheng, Yang Ai, Zhen-Hua Ling</i>	

Consonant-Emphasis Method Incorporating Robust Consonant-Section Detection to Improve Intelligibility of Bone-Conducted Speech	849
<i>Yasufumi Uezu, Sicheng Wang, Teruki Toya, Masashi Unoki</i>	
Downstream Task Agnostic Speech Enhancement with Self-Supervised Representation Loss	854
<i>Hiroshi Sato, Ryo Masumura, Tsubasa Ochiai, Marc Delcroix, Takafumi Moriya, Takanori Ashihara, Kentaro Shinayama, Saki Mizuno, Mana Ihori, Tomohiro Tanaka, Nobukatsu Hojo</i>	
Perceptual Improvement of Deep Neural Network (DNN) Speech Coder Using Parametric and Non-Parametric Density Models	859
<i>Joon Byun, Seungmin Shin, Jongmo Sung, Seungkwon Beack, Youngcheol Park</i>	
DeFT-AN RT: Real-Time Multichannel Speech Enhancement Using Dense Frequency-Time Attentive Network and Non-Overlapping Synthesis Window	864
<i>Dongheon Lee, Dayun Choi, Jung-Woo Choi</i>	

SPEECH RECOGNITION: SIGNAL PROCESSING, ACOUSTIC MODELING, ROBUSTNESS, ADAPTATION 2

A More Accurate Internal Language Model Score Estimation for the Hybrid Autoregressive Transducer	869
<i>Kyungmin Lee, Haeri Kim, Sichen Jin, Jinhwan Park, Youngho Han</i>	
Attention Gate Between Capsules in Fully Capsule-Network Speech Recognition	874
<i>Kyungmin Lee, Hyeontaek Lim, Mun-Hwan Lee, Hong-Gee Kim</i>	
ML-SUPERB: Multilingual Speech Universal PERFORMANCE Benchmark	884
<i>Jiatong Shi, Dan Berrebbi, William Chen, En-Pei Hu, Wei-Ping Huang, Ho-Lam Chung, Xuankai Chang, Shang-Wen Li, Abdelrahman Mohamed, Hung-Yi Lee, Shinji Watanabe</i>	
General-Purpose Adversarial Training for Enhanced Automatic Speech Recognition Model Generalization	889
<i>Dohee Kim, Daeyeol Shim, Joon-Hyuk Chang</i>	
Joint Instance Reconstruction and Feature Subspace Alignment for Cross-Domain Speech Emotion Recognition	894
<i>Keke Zhao, Peng Song, Shaokai Li, Wenming Zheng</i>	
Knowledge Distillation for Neural Transducer-Based Target-Speaker ASR: Exploiting Parallel Mixture/Single-Talker Speech Data	899
<i>Takafumi Moriya, Hiroshi Sato, Tsubasa Ochiai, Marc Delcroix, Takanori Ashihara, Kohei Matsuura, Tomohiro Tanaka, Ryo Masumura, Atsunori Ogawa, Taichi Asami</i>	
Random Utterance Concatenation Based Data Augmentation for Improving Short-Video Speech Recognition	904
<i>Yist Y. Lin, Tao Han, Haihua Xu, Van Tung Pham, Yerbolat Khassanov, Tze Yuang Chong, Yi He, Lu Lu, Zejun Ma</i>	
Adapter Incremental Continual Learning of Efficient Audio Spectrogram Transformers	909
<i>Nithish Muthuchamy Selvaraj, Xiaobao Guo, Adams Kong, Bingquan Shen, Alex Kot</i>	
Rethinking Speech Recognition with a Multimodal Perspective Via Acoustic and Semantic Cooperative Decoding	914
<i>Tian-Hao Zhang, Hai-Bo Qin, Zhi-Hao Lai, Song-Lu Chen, Qi Liu, Feng Chen, Xinyuan Qian, Xu-Cheng Yin</i>	

Improving Code-Switching and Name Entity Recognition in ASR with Speech Editing Based Data Augmentation	919
<i>Zheng Liang, Zheshu Song, Ziyang Ma, Chenpeng Du, Kai Yu, Xie Chen</i>	
Bypass Temporal Classification: Weakly Supervised Automatic Speech Recognition with Imperfect Transcripts	924
<i>Dongji Gao, Matthew Wiesner, Hainan Xu, Leibny Paola Garcia, Daniel Povey, Sanjeev Khudanpur</i>	
DCCRN-KWS: An Audio Bias Based Model for Noise Robust Small-Footprint Keyword Spotting.....	929
<i>Shubo Lv, Xiong Wang, Sining Sun, Long Ma, Lei Xie</i>	
OTF: Optimal Transport Based Fusion of Supervised and Self-Supervised Learning Models for Automatic Speech Recognition	934
<i>Li Fu, Siqi Li, Qingtao Li, Fangzhu Li, Liping Deng, Lu Fan, Meng Chen, Youzheng Wu, Xiaodong He</i>	
Approximate Nearest Neighbour Phrase Mining for Contextual Speech Recognition.....	939
<i>Maurits Bleeker, Pawel Swietojanski, Stefan Braun, Xiaodan Zhuang</i>	
Rehearsal-Free Online Continual Learning for Automatic Speech Recognition.....	944
<i>Steven Vander Eeckt, Hugo Van Hamme</i>	

SPEECH RECOGNITION: TECHNOLOGIES AND SYSTEMS FOR NEW APPLICATIONS 2

Phonetic and Prosody-Aware Self-Supervised Learning Approach for Non-Native Fluency Scoring	949
<i>Kaiqi Fu, Shaojun Gao, Shuju Shi, Xiaohai Tian, Wei Li, Zejun Ma</i>	
Disentangling the Contribution of Non-Native Speech in Automated Pronunciation Assessment.....	954
<i>Shuju Shi, Kaiqi Fu, Yiwei Gu, Xiaohai Tian, Shaojun Gao, Wei Li, Zejun Ma</i>	
A Joint Model for Pronunciation Assessment and Mispronunciation Detection and Diagnosis with Multi-Task Learning	959
<i>Hyungshin Ryu, Sunhee Kim, Minhwa Chung</i>	
Assessing Intelligibility in Non-Native Speech: Comparing Measures Obtained at Different Levels	964
<i>Xing Wei, Roeland Van Hout, Catia Cucchiarini, Danielle Reuvekamp, Helmer Strik</i>	
End-To-End Word-Level Pronunciation Assessment with MASK Pre-Training.....	969
<i>Yukang Liang, Kaitao Song, Shaoguang Mao, Huiqiang Jiang, Luna Qiu, Yuqing Yang, Dongsheng Li, Linli Xu, Lili Qiu</i>	
A Hierarchical Context-Aware Modeling Approach for Multi-Aspect and Multi-Granular Pronunciation Assessment	974
<i>Fu-An Chao, Tien-Hong Lo, Tzu-I Wu, Yao-Ting Sung, Berlin Chen</i>	
Automatic Prediction of Language Learners' Listenability Using Speech and Text Features Extracted from Listening Drills.....	979
<i>Yingxiang Gao, Jaehyun Choi, Nobuaki Minematsu, Noriko Nakanishi, Daisuke Saito</i>	
Assessment of Non-Native Speech Intelligibility Using Wav2vec2-Based Mispronunciation Detection and Multi-Level Goodness of Pronunciation Transformer.....	984
<i>Ram C. M. C. Shekar, Mu Yang, Kevin Hirschi, Stephen Looney, Okim Kang, John H. L. Hansen</i>	

Adapting an Unadaptable ASR System	989
<i>Rao Ma, Mengjie Qian, Mark J. F. Gales, Kate M. Knill</i>	
Addressing Cold Start Problem for End-To-End Automatic Speech Scoring.....	994
<i>Jungbae Park, Seungtaek Choi</i>	
Improving Grapheme-To-Phoneme Conversion by Learning Pronunciations from Speech Recordings.....	999
<i>Sam Ribeiro, Giulia Comini, Jaime Lorenzo-Trueba</i>	
Orthography-Based Pronunciation Scoring for Better CAPT Feedback	1004
<i>Caitlin Richter, Ragnar Pálsson, Luke O'Brien, Kolbrún Friðriksdóttir, Branislav Bédi, Eydis Huld Magnúsdóttir, Jón Guðnason</i>	
Zero-Shot Automatic Pronunciation Assessment	1009
<i>Hongfu Liu, Mingqian Shi, Ye Wang</i>	
Mispronunciation Detection and Diagnosis Model for Tonal Language, Applied to Vietnamese.....	1014
<i>Tuong Tu Huu, Viet Thanh Pham, Thi Thu Trang Nguyen, Thai Lai Dao</i>	

KEYNOTE 2

Beyond the AI Hype: Balancing Innovation and Social Responsibility	1019
<i>Virginia Dignum</i>	

PARALINGUISTICS 1

Detection of Emotional Hotspots in Meetings Using a Cross-Corpus Approach	1020
<i>Georg Stemmer, Paulo Lopez Meyer, Juan Del Hoyo Ontiveros, Jose Lopez, Hector A. Cordourier, Tobias Bocklet</i>	
Detection of Laughter and Screaming Using the Attention and CTC Models.....	1025
<i>Takuto Matsuda, Yoshiko Arimoto</i>	
Capturing Formality in Speech Across Domains and Languages.....	1030
<i>Debasmita Bhattacharya, Jie Chi, Julia Hirschberg, Peter Bell</i>	
Towards Robust Family-Infant Audio Analysis Based on Unsupervised Pretraining of Wav2vec 2.0 on Large-Scale Unlabeled Family Audio.....	1035
<i>Jialu Li, Mark Hasegawa-Johnson, Nancy L. McElwain</i>	
Cues to Next-Speaker Projection in Conversational Swedish: Evidence from Reaction Times.....	1040
<i>Kathrin Feindt, Martina Rossi, Ghazaleh Esfandiari-Baiat, Axel G. Ekström, Margaret Zellers</i>	
Multiple Instance Learning for Inference of Child Attachment from Paralinguistic Aspects of Speech	1045
<i>Areej Buker, Huda Alsofyani, Alessandro Vinciarelli</i>	

SPEECH ENHANCEMENT AND DENOISING

Real-Time Joint Personalized Speech Enhancement and Acoustic Echo Cancellation	1050
<i>Sefik Emre Eskimez, Takuya Yoshioka, Alex Ju, Min Tang, Tanel Pärnamaa, Huaming Wang</i>	

TaylorBeamixer: Learning Taylor-Inspired All-Neural Multi-Channel Speech Enhancement from Beam-Space Dictionary Perspective.....	1055
<i>Andong Li, Weixin Meng, Guochen Yu, Wenzhe Liu, Xiaodong Li, Chengshi Zheng</i>	
MFT-CRN:Multi-Scale Fourier Transform for Monaural Speech Enhancement	1060
<i>Yulong Wang, Xueliang Zhang</i>	
Variance-Preserving-Based Interpolation Diffusion Models for Speech Enhancement	1065
<i>Zilu Guo, Jun Du, Chin-Hui Lee, Yu Gao, Wenbin Zhang</i>	
Multi-Input Multi-Output Complex Spectral Mapping for Speaker Separation.....	1070
<i>Hassan Taherian, Ashutosh Pandey, Daniel Wong, Buye Xu, Deliang Wang</i>	
Short-Term Extrapolation of Speech Signals Using Recursive Neural Networks in the STFT Domain	1075
<i>Maurice Oberhag, Daniel Neudek, Rainer Martin, Tobias Rosenkranz, Henning Puder</i>	

SPEECH SYNTHESIS: EVALUATION

Listener Sensitivity to Deviating Obstruents in WaveNet	1080
<i>Ayushi Pandey, Jens Edlund, Sébastien Le Maguer, Naomi Harte</i>	
How Generative Spoken Language Modeling Encodes Noisy Speech: Investigation from Phonetics to Syntactics	1085
<i>Joonyong Park, Shinnosuke Takamichi, Tomohiko Nakamura, Kentaro Seki, Detai Xin, Hiroshi Saruwatari</i>	
MOS vs. AB: Evaluating Text-To-Speech Systems Reliably Using Clustered Standard Errors.....	1090
<i>Joshua Camp, Tom Kenter, Lev Finkelstein, Rob Clark</i>	
RAMP: Retrieval-Augmented MOS Prediction Via Confidence-Based Dynamic Weighting.....	1095
<i>Hui Wang, Shiwan Zhao, Xiguang Zheng, Yong Qin</i>	
Can Better Perception Become a Disadvantage? Synthetic Speech Perception in Congenitally Blind Users.....	1100
<i>Gerda Ana Melnik-Leroy, Gediminas Navickas</i>	
Investigating Range-Equalizing Bias in Mean Opinion Score Ratings of Synthesized Speech	1104
<i>Erica Cooper, Junichi Yamagishi</i>	

END-TO-END SPOKEN DIALOG SYSTEMS

Can ChatGPT Detect Intent? Evaluating Large Language Models for Spoken Language Understanding	1109
<i>Mutian He, Philip N. Garner</i>	
Improving End-To-End SLU Performance with Prosodic Attention and Distillation.....	1114
<i>Shangeth Rajaa</i>	
Modality Confidence Aware Training for Robust End-To-End Spoken Language Understanding	1119
<i>Suyoun Kim, Akshat Shrivastava, Duc Le, Ju Lin, Ozlem Kalinli, Michael L. Seltzer</i>	
Cross-Modal Semantic Alignment Before Fusion for Two-Pass End-To-End Spoken Language Understanding	1124
<i>Lingyan Huang, Tao Li, Haodong Zhou, Qingyang Hong, Lin Li</i>	

ConvKT: Conversation-Level Knowledge Transfer for Context Aware End-To-End Spoken Language Understanding1129
Vishal Sunder, Eric Fosler-Lussier, Samuel Thomas, Hong-Kwang J Kuo, Brian Kingsbury

GhostT5: Generate More Features with Cheap Operations to Improve Textless Spoken Question Answering1134
Xuxin Cheng, Zhihong Zhu, Ziyu Yao, Hongxiang Li, Yaowei Li, Yuexian Zou

BIOSIGNAL-ENABLED SPOKEN COMMUNICATION

Obstructive Sleep Apnea Detection Using Pre-Trained Speech Representations1139
Kaibo Zhang, Lili Cao, Yiming Ding, Yanru Li, Chao Zhang, Ji Wu, Demin Han

EEG-Based Auditory Attention Detection with Spatiotemporal Graph and Graph Convolutional Network1144
Ruicong Wang, Siqi Cai, Haizhou Li

Silent Speech Recognition with Articulator Positions Estimated from Tongue Ultrasound and Lip Video1149
Rachel Beeson, Korin Richmond

Auditory Attention Detection in Real-Life Scenarios Using Common Spatial Patterns from EEG1154
Kai Yang, Zhuang Xie, Di Zhou, Longbiao Wang, Gaoyan Zhang

Diff-E: Diffusion-Based Learning for Decoding Imagined Speech EEG1159
Soowon Kim, Young-Eun Lee, Seo-Hyun Lee, Seong-Whan Lee

Towards Ultrasound Tongue Image Prediction from EEG During Speech Production1164
Tamás Gábor Csapó, Frigyes Viktor Arthur, Péter Nagy, Adám Boncz

Adaptation of Tongue Ultrasound-Based Silent Speech Interfaces Using Spatial Transformer Networks1169
László Tóth, Amin Honarmandi Shandiz, Gábor Gosztolya, Tamás Gábor Csapó

STE-GAN: Speech-To-Electromyography Signal Conversion Using Generative Adversarial Networks1174
Kevin Scheck, Tanja Schultz

Spanish Phone Confusion Analysis for EMG-Based Silent Speech Interfaces1179
Inge Salomons, Eder Del Blanco, Eva Navas, Inma Hernández

Hybrid Silent Speech Interface Through Fusion of Electroencephalography and Electromyography1184
Huiyan Li, Mingyi Wang, Han Gao, Shuo Zhao, Guang Li, You Wang

NEURAL-BASED SPEECH AND ACOUSTIC ANALYSIS

Can Self-Supervised Neural Representations Pre-Trained on Human Speech Distinguish Animal Callers?1189
Eklavya Sarkar, Mathew Magimai.-Doss

Discovering COVID-19 Coughing and Breathing Patterns from Unlabeled Data Using Contrastive Learning with Varying Pre-Training Domains1194
Jinjin Cai, Sudip Vhaduri, Xiao Luo

Background-Aware Modeling for Weakly Supervised Sound Event Detection	1199
<i>Yifei Xin, Dongchao Yang, Yuexian Zou</i>	
How to (Virtually) Train Your Speaker Localizer.....	1204
<i>Prerak Srivastava, Antoine Deleforge, Archontis Politis, Emmanuel Vincent</i>	
MMER: Multimodal Multi-Task Learning for Speech Emotion Recognition.....	1209
<i>Sreyan Ghosh, Utkarsh Tyagi, S Ramaneswaran, Harshvardhan Srivastava, Dinesh Manocha</i>	
A Multi-Task Learning Framework for Sound Event Detection Using High-Level Acoustic Characteristics of Sounds	1214
<i>Tanmay Khandelwal, Rohan Kumar Das</i>	

DIGO - DIALOG FOR GOOD: SPEECH AND LANGUAGE TECHNOLOGY FOR SOCIAL GOOD

A Multimodal Investigation of Speech, Text, Cognitive and Facial Video Features for Characterizing Depression with and Without Medication.....	1219
<i>Michael Neumann, Hardik Kothare, Doug Habberstad, Vikram Ramanarayanan</i>	
Understanding Disrupted Sentences Using Underspecified Abstract Meaning Representation	1224
<i>Angus Addlesee, Marco Damonte</i>	
Developing Speech Processing Pipelines for Police Accountability	1229
<i>Anjalie Field, Prateek Verma, Nay San, Jennifer L. Eberhardt, Dan Jurafsky</i>	
Prosody-Controllable Gender-Ambiguous Speech Synthesis: A Tool for Investigating Implicit Bias in Speech Perception	1234
<i>Éva Székely, Joakim Gustafson, Ilaria Torre</i>	
Affective Attributes of French Caregivers' Professional Speech.....	1239
<i>Jean-Luc Rouas, Yaru Wu, Takaaki Shochi</i>	

SPEECH RECOGNITION: SIGNAL PROCESSING, ACOUSTIC MODELING, ROBUSTNESS, ADAPTATION 3

ASR Data Augmentation in Low-Resource Settings Using Cross-Lingual Multi-Speaker TTS and Cross-Lingual Voice Conversion.....	1244
<i>Edresson Casanova, Christopher Shulby, Alexander Korolev, Arnaldo Candido Junior, Anderson Da Silva Soares, Sandra Aluísio, Moacir Antonelli Ponti</i>	
Personality-Aware Training Based Speaker Adaptation for End-To-End Speech Recognition	1249
<i>Yue Gu, Zhihao Du, Shiliang Zhang, Qian Chen, Jiqing Han</i>	
Target Vocabulary Recognition Based on Multi-Task Learning with Decomposed Teacher Sequences	1254
<i>Aoi Ito, Tatsuya Komatsu, Yusuke Fujita, Yusuke Kida</i>	
Wave to Syntax: Probing Spoken Language Models for Syntax	1259
<i>Gaofei Shen, Afra Alishahi, Arianna Bisazza, Grzegorz Chrupala</i>	
Effective Training of Attention-Based Contextual Biasing Adapters with Synthetic Audio for Personalised ASR	1264
<i>Burin Naowarat, Philip Harding, Pasquale D'Alterio, Sibong Tong, Bashar Awwad Shiekh Hasan</i>	

Pushing the Limits of Unsupervised Unit Discovery for SSL Speech Representation.....	1269
<i>Ziyang Ma, Zhisheng Zheng, Guanrou Yang, Yu Wang, Chao Zhang, Xie Chen</i>	
SlothSpeech: Denial-Of-Service Attack Against Speech Recognition Models	1274
<i>Mirazul Haque, Rutvij Shah, Simin Chen, Berrak Sisman, Cong Liu, Wei Yang</i>	
CLRL-Tuning: A Novel Continual Learning Approach for Automatic Speech Recognition	1279
<i>Zhihan Wang, Feng Hou, Ruili Wang</i>	
Exploring Sources of Racial Bias in Automatic Speech Recognition Through the Lens of Rhythmic Variation	1284
<i>Li-Fang Lai, Nicole Holliday</i>	
Can Contextual Biasing Remain Effective with Whisper and GPT-2?.....	1289
<i>Guangzhi Sun, Xianrui Zheng, Chao Zhang, Philip C. Woodland</i>	
Masked Modeling Duo for Speech: Specializing General-Purpose Audio Representation to Speech Using Denoising Distillation	1294
<i>Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, Kunio Kashino</i>	
Improving RNN Transducer Acoustic Models for English Conversational Speech Recognition.....	1299
<i>Xiaodong Cui, George Saon, Brian Kingsbury</i>	
MixRep: Hidden Representation Mixup for Low-Resource Speech Recognition.....	1304
<i>Jiamin Xie, John H. L. Hansen</i>	
Adapting Multi-Lingual ASR Models for Handling Multiple Talkers.....	1314
<i>Chenda Li, Yao Qian, Zhuo Chen, Naoyuki Kanda, Dongmei Wang, Takuya Yoshioka, Yanmin Qian, Michael Zeng</i>	
Adapter-Tuning with Effective Token-Dependent Representation Shift for Automatic Speech Recognition	1319
<i>Dianwen Ng, Chong Zhang, Ruixi Zhang, Yukun Ma, Trung Hieu Nguyen, Chongjia Ni, Shengkui Zhao, Qian Chen, Wen Wang, Eng Siong Chng, Bin Ma</i>	
Model-Internal Slot-Triggered Biasing for Domain Expansion in Neural Transducer ASR Models	1324
<i>Yiting Lu, Philip Harding, Kanthashree Mysore Sathyendra, Sibor Tong, Xuandi Fu, Jing Liu, Feng-Ju Chang, Simon Wiesler, Grant P. Strimel</i>	
Delay-Penalized CTC Implemented Based on Finite State Transducer	1329
<i>Zengwei Yao, Wei Kang, Fangjun Kuang, Liyong Guo, Xiaoyu Yang, Yifan Yang, Long Lin, Daniel Povey</i>	

SPEECH RECOGNITION: ARCHITECTURE, SEARCH, AND LINGUISTIC COMPONENTS

2

Text-Only Domain Adaptation for End-To-End Speech Recognition Through Down-Sampling Acoustic Representation	1334
<i>Jiaxu Zhu, Weinan Tong, Yaoxun Xu, Changhe Song, Zhiyong Wu, Zhao You, Dan Su, Dong Yu, Helen Meng</i>	
Knowledge Distillation Approach for Efficient Internal Language Model Estimation	1339
<i>Zhipeng Chen, Haihua Xu, Yerbolat Khassanov, Yi He, Lu Lu, Zejun Ma, Ji Wu</i>	
Language Model Personalization for Improved Touchscreen Typing	1344
<i>Jiban Adhikary, Keith Vertanen</i>	

Blank Collapse: Compressing CTC Emission for the Faster Decoding	1349
<i>Minkyu Jung, Ohhyeok Kwon, Seunghyun Seo, Soonshin Seo</i>	
Improving Joint Speech-Text Representations Without Alignment.....	1354
<i>Cal Peyser, Zhong Meng, Rohit Prabhavalkar, Andrew Rosenberg, Tara Sainath, Michael Picheny, Kyunghyun Cho, Ke Hu</i>	
Leveraging Cross-Utterance Context for ASR Decoding.....	1359
<i>Robert Flynn, Anton Ragni</i>	
Knowledge Transfer from Pre-Trained Language Models to Cif-Based Speech Recognizers Via Hierarchical Distillation	1364
<i>Minglun Han, Feilong Chen, Jing Shi, Shuang Xu, Bo Xu</i>	
Integration of Frame- And Label-Synchronous Beam Search for Streaming Encoder-Decoder Speech Recognition.....	1369
<i>Emiru Tsunoo, Hayato Futami, Yosuke Kashiwagi, Siddhant Arora, Shinji Watanabe</i>	
A Neural Time Alignment Module for End-To-End Automatic Speech Recognition.....	1374
<i>Dongcheng Jiang, Chao Zhang, Philip C. Woodland</i>	
Accelerating Transducers Through Adjacent Token Merging	1379
<i>Yuang Li, Yu Wu, Jinyu Li, Shujie Liu</i>	
Language-Universal Phonetic Representation in Multilingual Speech Pretraining for Low-Resource Speech Recognition.....	1384
<i>Siyuan Feng, Ming Tu, Rui Xia, Chuanzeng Huang, Yuxuan Wang</i>	
Language-Routing Mixture of Experts for Multilingual and Code-Switching Speech Recognition	1389
<i>Wenxuan Wang, Guodong Ma, Yuke Li, Binbin Du</i>	
Embedding Articulatory Constraints for Low-Resource Speech Recognition Based on Large Pre-Trained Model	1394
<i>Jaeyoung Lee, Masato Mimura, Tatsuya Kawahara</i>	

VOLUME 3

Exploration of Efficient End-To-End ASR Using Discretized Input from Self-Supervised Learning.....	1399
<i>Xuankai Chang, Brian Yan, Yuya Fujita, Takashi Maekaku, Shinji Watanabe</i>	
SpellMapper: A Non-Autoregressive Neural Spellchecker for ASR Customization with Candidate Retrieval Based on N-Gram Mappings.....	1404
<i>Alexandra Antonova, Evelina Bakhturina, Boris Ginsburg</i>	
Text Injection for Capitalization and Turn-Taking Prediction in Speech Models.....	1409
<i>Shaan Bijwadia, Shuo-Yiin Chang, Weiran Wang, Zhong Meng, Hao Zhang</i>	
Confidence-Based Ensembles of End-To-End Speech Recognition Models.....	1414
<i>Igor Gitman, Vitaly Lavrukhin, Aleksandr Laptev, Boris Ginsburg</i>	
Unsupervised Code-Switched Text Generation from Parallel Text	1419
<i>Jie Chi, Brian Lu, Jason Eisner, Peter Bell, Preethi Jyothi, Ahmed M. Ali</i>	
A Binary Keyword Spotting System with Error-Diffusion Based Feature Binarization.....	1424
<i>Dingyi Wang, Mengjie Luo, Lin Li, Xiaoqin Wang, Shushan Qiao, Yumei Zhou</i>	

Language-Universal Phonetic Encoder for Low-Resource Speech Recognition	1429
<i>Siyuan Feng, Ming Tu, Rui Xia, Chuanzeng Huang, Yuxuan Wang</i>	
A Lexical-Aware Non-Autoregressive Transformer-Based ASR Model	1434
<i>Chong-En Lin, Kuan-Yu Chen</i>	
Improving Under-Resourced Code-Switched Speech Recognition: Large Pre-Trained Models Or Architectural Interventions	1439
<i>Joshua Jansen Van Vuren, Thomas Niesler</i>	

SPOKEN LANGUAGE TRANSLATION, INFORMATION RETRIEVAL, SUMMARIZATION, RESOURCES, AND EVALUATION 1

Pragmatic Pertinence: A Learnable Confidence Metric to Assess the Subjective Quality of LM-Generated Text.....	1444
<i>Jerome R. Bellegarda</i>	
ASR and Emotional Speech: A Word-Level Investigation of the Mutual Impact of Speech and Emotion Recognition.....	1449
<i>Yuanchao Li, Zeyu Zhao, Ondrej Klejch, Peter Bell, Catherine Lai</i>	
BASS: Block-Wise Adaptation for Speech Summarization	1454
<i>Roshan Sharma, Siddhant Arora, Kenneth Zheng, Shinji Watanabe, Rita Singh, Bhiksha Raj</i>	
Speaker Tracking Using Graph Attention Networks with Varying Duration Utterances Across Multi-Channel Naturalistic Data: Fearless Steps Apollo-11 Audio Corpus.....	1459
<i>Meena M. C. Shekar, John H. L. Hansen</i>	
Combining Language Corpora in a Japanese Electromagnetic Articulography Database for Acoustic-To-Articulatory Inversion.....	1464
<i>Tianfang Yan, Kikuo Maekawa, Yukiko Nota, Masayuki Hirata</i>	
A Dual Attention-Based Modality-Collaborative Fusion Network for Emotion Recognition.....	1468
<i>Xiaoheng Zhang, Yang Li</i>	
Large Dataset Generation of Synchronized Music Audio and Lyrics at Scale Using Teacher-Student Paradigm.....	1473
<i>Cristian Chivriga, Rinita Roy</i>	
Enc-Dec RNN Acoustic Word Embeddings Learned Via Pairwise Prediction.....	1478
<i>Adhiraj Banerjee, Vipul Arora</i>	
Query Based Acoustic Summarization for Podcasts.....	1483
<i>Samantha Kotey, Rozenn Dahyot, Naomi Harte</i>	
Spot Keywords from Very Noisy and Mixed Speech	1488
<i>Ying Shi, Dong Wang, Lantian Li, Jiqing Han, Shi Yin</i>	
Knowledge Distillation on Joint Task End-To-End Speech Translation.....	1493
<i>Khandokar Md. Nayem, Ran Xue, Ching-Yun Chang, Akshaya Vishnu Kudlu Shanbhogue</i>	
Investigating Pre-Trained Audio Encoders in the Low-Resource Condition.....	1498
<i>Hao Yang, Jinming Zhao, Gholamreza Haffari, Ehsan Shareghi</i>	
Improving Textless Spoken Language Understanding with Discrete Units as Intermediate Target	1503
<i>Guan-Wei Wu, Guan-Ting Lin, Shang-Wen Li, Hung-Yi Lee</i>	

SPEECH, VOICE, AND HEARING DISORDERS 1

Debiased Automatic Speech Recognition for Dysarthric Speech Via Sample Reweighting with Sample Affinity Test.....	1508
<i>Eungbeom Kim, Yunkee Chae, Jaeheon Sim, Kyogu Lee</i>	
Multimodal Locally Enhanced Transformer for Continuous Sign Language Recognition	1513
<i>Katerina Papadimitriou, Gerasimos Potamianos</i>	
Towards Supporting an Early Diagnosis of Multiple Sclerosis Using Vocal Features	1518
<i>Monica Gonzalez-Machorro, Pascal Hecker, Uwe D. Reichel, Helly N. Hammer, Robert Hoepner, Lisa Pedrotti, Alisha Zmutt, Hesam Sagha, Johan Van Beek, Florian Eyben, Dagmar M. Schuller, Björn W. Schuller, Bert Arnrich</i>	
Whisper Features for Dysarthric Severity-Level Classification	1523
<i>Siddharth Rathod, Monil Charola, Akshat Vora, Yash Jogi, Hemant A. Patil</i>	
A New Benchmark of Aphasia Speech Recognition and Detection Based on E-Branchformer and Multi-Task Learning	1528
<i>Jiyang Tang, William Chen, Xuankai Chang, Shinji Watanabe, Brian Macwhinney</i>	
Dysarthric Speech Recognition, Detection and Classification Using Raw Phase and Magnitude Spectra.....	1533
<i>Zhengjun Yue, Erfan Loweimi, Zoran Cvetkovic</i>	
A Stutter Seldom Comes Alone – Cross-Corpus Stuttering Detection as a Multi-Label Problem	1538
<i>Sebastian P. Bayerl, Dominik Wagner, Ilja Baumann, Florian Hönig, Tobias Bocklet, Elmar Nöth, Korbinian Riedhammer</i>	
Transfer Learning to Aid Dysarthria Severity Classification for Patients with Amyotrophic Lateral Sclerosis	1543
<i>Tanuka Bhattacharjee, Anjali Jayakumar, Yamini Belur, Atchayaram Nalini, Ravi Yadav, Prasanta Kumar Ghosh</i>	
DuTa-VC: A Duration-Aware Typical-To-Atypical Voice Conversion Approach with Diffusion Probabilistic Model	1548
<i>Helin Wang, Thomas Thebaud, Jesús Villalba, Myra Sydnor, Becky Lammers, Najim Dehak, Laureano Moro-Velazquez</i>	
CNVVE: Dataset and Benchmark for Classifying Non-Verbal Voice	1553
<i>Ramin Hedeshy, Raphael Menges, Steffen Staab</i>	
Arabic Dysarthric Speech Recognition Using Adversarial and Signal-Based Augmentation	1558
<i>Massa Baali, Ibrahim Almakky, Shady Shehata, Fakhri Karray</i>	
Weakly-Supervised Forced Alignment of Disfluent Speech Using Phoneme-Level Modeling.....	1563
<i>Theodoros Kouzelis, Georgios Paraskevopoulos, Athanasios Katsamanis, Vassilis Katsouros</i>	
Glottal Source Analysis of Voice Deficits in Basal Ganglia Dysfunction: Evidence from De Novo Parkinson's Disease and Huntington's Disease.....	1568
<i>Michal Novotný, Tereza Tykalová, Michal Šimek, Tomáš Kouba, Jan Ruz</i>	
An Analysis of Glottal Features of Chronic Kidney Disease Speech and Its Application to CKD Detection	1573
<i>Jihyun Mun, Sunhee Kim, Myeong Ju Kim, Jiwon Ryu, Sejoong Kim, Minhwa Chung</i>	

Weakly Supervised Glottis Segmentation in High-Speed Videoendoscopy Using Bounding Box Labels	1578
<i>Varun Belagali, Achuth Rao, Prasanta Kumar Ghosh</i>	

SPEECH RECOGNITION: TECHNOLOGIES AND SYSTEMS FOR NEW APPLICATIONS 3

An Efficient and Noise-Robust Audiovisual Encoder for Audiovisual Speech Recognition.....	1583
<i>Zhengyang Li, Chenwei Liang, Timo Lohrenz, Marvin Sach, Björn Möller, Tim Fingscheidt</i>	
A Novel Self-Training Approach for Low-Resource Speech Recognition	1588
<i>Satwinder Singh, Feng Hou, Ruili Wang</i>	
FunASR: A Fundamental End-To-End Speech Recognition Toolkit.....	1593
<i>Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Shiliang Zhang</i>	
Streaming Audio-Visual Speech Recognition with Alignment Regularization	1598
<i>Pingchuan Ma, Niko Moritz, Stavros Petridis, Christian Fuegen, Maja Pantic</i>	
SparseVSR: Lightweight and Noise Robust Visual Speech Recognition	1603
<i>Adriana Fernandez-Lopez, Honglie Chen, Pingchuan Ma, Alexandros Haliassos, Stavros Petridis, Maja Pantic</i>	
Multimodal Speech Recognition for Language-Guided Embodied Agents.....	1608
<i>Allen Chang, Xiaoyuan Zhu, Aarav Monga, Seoho Ahn, Tejas Srinivasan, Jesse Thomason</i>	

SPOKEN TERM DETECTION AND VOICE SEARCH

Matching Latent Encoding for Audio-Text Based Keyword Spotting.....	1613
<i>Kumari Nishu, Minsik Cho, Devang Naik</i>	
Self-Paced Pattern Augmentation for Spoken Term Detection in Zero-Resource	1618
<i>Sudhakar P, Sreenivasa K. Rao, Pabitra Mitra</i>	
On-Device Constrained Self-Supervised Speech Representation Learning for Keyword Spotting Via Knowledge Distillation	1623
<i>Gene-Ping Yang, Yue Gu, Qingming Tang, Dongsu Du, Yuzong Liu</i>	
Online Continual Learning in Keyword Spotting for Low-Resource Devices Via Pooling High-Order Temporal Statistics	1628
<i>Umberto Michieli, Pablo Peso Parada, Mete Ozay</i>	
Improving Small Footprint Few-Shot Keyword Spotting with Supervision on Auxiliary Data.....	1633
<i>Seunghan Yang, Byeonggeun Kim, Kyuhong Shim, Simyoung Chang</i>	
Robust Keyword Spotting for Noisy Environments by Leveraging Speech Enhancement and Speech Presence Probability.....	1638
<i>Chouchang Yang, Yashas Malur Saidutta, Rakshith Sharma Srinivasa, Ching-Hua Lee, Yilin Shen, Hongxia Jin</i>	

MODELS FOR STREAMING ASR

Enhancing the Unified Streaming and Non-Streaming Model with Contrastive Learning	1643
<i>Yuting Yang, Yuke Li, Binbin Du</i>	

ZeroPrompt: Streaming Acoustic Encoders Are Zero-Shot Masked LMs.....	1648
<i>Xingchen Song, Di Wu, Binbin Zhang, Zhendong Peng, Bo Dang, Fuping Pan, Zhiyong Wu</i>	
Improved Training for End-To-End Streaming Automatic Speech Recognition Model with Punctuation.....	1653
<i>Hanbyul Kim, Seunghyun Seo, Lukas Lee, Seolki Baek</i>	
DCTX-Conformer: Dynamic Context Carry-Over for Low Latency Unified Streaming and Non-Streaming Conformer.....	1658
<i>Goeric Huybrechts, Srikanth Ronanki, Xilai Li, Hadis Nosrati, Sravan Bodapati, Katrin Kirchhoff</i>	
Knowledge Distillation from Non-Streaming to Streaming ASR Encoder Using Auxiliary Non-Streaming Layer.....	1663
<i>Kyuhong Shim, Jinkyu Lee, Simyoung Chang, Kyuwoong Hwang</i>	
Adaptive Contextual Biasing for Transducer Based Streaming Speech Recognition.....	1668
<i>Tianyi Xu, Zhanheng Yang, Kaixun Huang, Pengcheng Guo, Ao Zhang, Biao Li, Changru Chen, Chao Li, Lei Xie</i>	

SOURCE SEPARATION

Audio-Visual Speech Separation in Noisy Environments with a Lightweight Iterative Model.....	1673
<i>Héctor Martel, Julius Richter, Kai Li, Xiaolin Hu, Timo Gerkmann</i>	
Remixing-Based Unsupervised Source Separation from Scratch.....	1678
<i>Kohei Saijo, Tetsuji Ogawa</i>	
CAPTDURE: Captioned Sound Dataset of Single Sources.....	1683
<i>Yuki Okamoto, Kanta Shimonishi, Keisuke Imoto, Kota Dohi, Shota Horiguchi, Yohei Kawaguchi</i>	
Recursive Sound Source Separation with Deep Learning-Based Beamforming for Unknown Number of Sources.....	1688
<i>Hokuto Munakata, Ryu Takeda, Kazunori Komatani</i>	
Multi-Channel Speech Separation with Cross-Attention and Beamforming.....	1693
<i>Ladislav Mosner, Oldrich Plchot, Junyi Peng, Lukáš Burget, Jan "Honza" Cernocký</i>	
Background-Sound Controllable Voice Source Separation.....	1698
<i>Deokjun Eom, Woo Hyun Nam, Kyung-Rae Kim</i>	

SPEECH AND LANGUAGE IN HEALTH: FROM REMOTE MONITORING TO MEDICAL CONVERSATIONS 1

An Automatic Multimodal Approach to Analyze Linguistic and Acoustic Cues on Parkinson's Disease Patients.....	1703
<i>Daniel Escobar-Grisales, Tomás Arias-Vergara, Cristian David Ríos-Urrego, Elmar Nöth, Adolfo M. García, Juan Rafael Orozco-Arroyave</i>	
Personalization for Robust Voice Pathology Detection in Sound Waves.....	1708
<i>Khanh-Tung Tran, Truong Hoang, Duy Khuong Nguyen, Hoang D. Nguyen, Xuan-Son Vu</i>	

Integrated and Enhanced Pipeline System to Support Spoken Language Analytics for Screening Neurocognitive Disorders.....	1713
<i>Helen Meng, Brian Mak, Man-Wai Mak, Helene Fung, Xianmin Gong, Timothy Kwok, Xunying Liu, Vincent Mok, Patrick Wong, Jean Woo, Xixin Wu, Ka Ho Wong, Shensheng Xu, Naijun Zheng, Ranzo Huang, Jiawen Kang, Xiaoquan Ke, Junan Li, Jinchao Li, Yi Wang</i>	
Capturing Mismatch Between Textual and Acoustic Emotion Expressions for Mood Identification in Bipolar Disorder.....	1718
<i>Minxue Niu, Amrit Romana, Mimansa Jaiswal, Melvin McInnis, Emily Mower Provost</i>	
FTA-Net: A Frequency and Time Attention Network for Speech Depression Detection.....	1723
<i>Qifei Li, Dong Wang, Yiming Ren, Yingming Gao, Ya Li</i>	
Bayesian Networks for the Robust and Unbiased Prediction of Depression and Its Symptoms Utilizing Speech and Multimodal Data.....	1728
<i>Salvatore Fara, Orlaith Hickey, Alexandra Georgescu, Stefano Gorla, Emilia Molimpakis, Nicholas Cummins</i>	
Hyper-Parameter Adaptation of Conformer ASR Systems for Elderly and Dysarthric Speech Recognition.....	1733
<i>Tianzi Wang, Shoukang Hu, Jiajun Deng, Zengrui Jin, Mengzhe Geng, Yi Wang, Helen Meng, Xunying Liu</i>	
Classifying Depression Symptom Severity: Assessment of Speech Representations in Personalized and Generalized Machine Learning Models.....	1738
<i>Edward L. Campbell, Judith Dineley, Pauline Conde, Faith Matcham, Katie M. White, Carolin Oetzmann, Sara Simblett, Stuart Bruce, Amos A. Folarin, Til Wykes, Srinivasan Vairavan, Richard J. B. Dobson, Laura Docio-Fernandez, Carmen Garcia-Mateo, Vaibhav A. Narayan, Matthew Hotopf, Nicholas Cummins</i>	
Active Learning for Abnormal Lung Sound Data Curation and Detection in Asthma.....	1743
<i>Shabnam Ghaffarzagdegan, Luca Bondi, Ho-Hsiang Wu, Sirajum Munir, Kelly J. Shields, Samarjit Das, Joseph Aracri</i>	
Automatic Assessment of Alzheimer's Across Three Languages Using Speech and Language Features.....	1748
<i>Paula A. Pérez-Toro, Tomás Arias-Vergara, Franziska Braun, Florian Hönig, Carlos A. Tobón-Quintero, David Aguillón, Francisco Lopera, Liliana Hincapié-Henao, Maria Schuster, Korbinian Riedhammer, Andreas Maier, Elmar Nöth, Juan Rafael Orozco-Arroyave</i>	
On-The-Fly Feature Based Rapid Speaker Adaptation for Dysarthric and Elderly Speech Recognition.....	1753
<i>Mengzhe Geng, Xurong Xie, Rongfeng Su, Jianwei Yu, Zengrui Jin, Tianzi Wang, Shujie Hu, Zi Ye, Helen Meng, Xunying Liu</i>	
Relationship Between LTAS-Based Spectral Moments and Acoustic Parameters of Hypokinetic Dysarthria in Parkinson's Disease.....	1758
<i>Jan Svihlik, Vojtech Illner, Petr Kryze, Mário Sousa, Paul Krack, Elina Tripoliti, Robert Jech, Jan Rusz</i>	
Respiratory Distress Estimation in Human-Robot Interaction Scenario.....	1763
<i>Eduardo Alvarado, Nicolás Grágeda, Alejandro Luzanto, Rodrigo Mahu, Jorge Wuth, Laura Mendoza, Richard Stern, Néstor Becerra Yoma</i>	

Prediction of the Gender-Based Violence Victim Condition Using Speech: What Do Machine Learning Models Rely On?	1768
<i>Emma Reyner-Fuentes, Esther Rituerto-González, Isabel Trancoso, Carmen Peláez-Moreno</i>	

Whisper Encoder Features for Infant Cry Classification.....	1773
<i>Monil Charola, Aastha Kachhi, Hemant A. Patil</i>	

SPEECH PERCEPTION

A Neural Architecture for Selective Attention to Speech Features	1778
<i>Nika Jurov, William Idsardi, Naomi H. Feldman</i>	

Quantifying Informational Masking Due to Masker Intelligibility in Same-Talker Speech-In-Speech Perception.....	1783
<i>Mingyue Huo, Yinglun Sun, Dan Fogerty, Yan Tang</i>	

On the Benefits of Self-Supervised Learned Speech Representations for Predicting Human Phonetic Misperceptions	1788
<i>Santiago Cuervo, Ricard Marxer</i>	

Predicting Perceptual Centers Located at Vowel Onset in German Speech Using Long Short-Term Memory Networks.....	1793
<i>Felicia Schulz, Mirella De Sisto, M. Paula Roncaglia-Denissen, Peter Hendrix</i>	

Exploring the Mutual Intelligibility Breakdown Caused by Sculpting Speech from a Competing Speech Signal	1798
<i>Martin Cooke, María Luisa García Lecumberri</i>	

Perception of Incomplete Voicing Neutralization of Obstruents in Tohoku Japanese	1803
<i>Mafuyu Kitahara, Naoya Watabe, Hiroto Noguchi, Chuyu Huang, Ayako Hashimoto, Ai Mizoguchi</i>	

PHONETICS AND PHONOLOGY: LANGUAGES AND VARIETIES

The Emergence of Obstruent-Intrinsic F0 and VOT as Cues to the Fortis/Lenis Contrast in West Central Bavarian.....	1808
<i>Jasmin Pöhnlein, Felicitas Kleber</i>	

(<i>l</i>) in Tsimane': A Preliminary Investigation	1813
<i>William N. Havard, Yaya Sy, Camila Scaff, Loann Peurey, Alejandrina Cristia</i>	

Segmental Features of Brazilian (Santa Catarina) Hunsrik	1818
<i>Dennis Hoffmann, Maria O'Reilly</i>	

Opening or Closing? An Electrolottographic Analysis of Voiceless Coda Consonants in Australian English.....	1823
<i>Louise Ratko, Joshua Penney, Felicity Cox</i>	

Increasing Aspiration of Word-Medial Fortis Plosives in Swiss Standard German.....	1828
<i>Franka Zebe</i>	

Lexical Stress and Velar Palatalization in Italian: A Spatio-Temporal Interaction	1833
<i>Bowei Shao, Philipp Buech, Anne Hermes, Maria Giavazzi</i>	

PARALINGUISTICS 2

Speaker Embeddings as Individuality Proxy for Voice Stress Detection.....	1838
<i>Zihan Wu, Neil Scheidwasser-Clow, Karl El Hajal, Milos Cernak</i>	
From Interval to Ordinal: A HMM Based Approach for Emotion Label Conversion.....	1843
<i>Jingyao Wu, Ting Dang, Vidhyasaharan Sethu, Eliathamby Ambikairajah</i>	
Turbo Your Multi-Modal Classification with Contrastive Learning.....	1848
<i>Zhiyu Zhang, Da Liu, Shengqiang Liu, Anna Wang, Jie Gao, Yali Li</i>	
Towards Paralinguistic-Only Speech Representations for End-To-End Speech Emotion Recognition.....	1853
<i>Georgios Ioannides, Michael Owen, Andrew Fletcher, Viktor Rozgic, Chao Wang</i>	
SOT: Self-Supervised Learning-Assisted Optimal Transport for Unsupervised Adaptive Speech Emotion Recognition.....	1858
<i>Ruiteng Zhang, Jianguo Wei, Xugang Lu, Yongwei Li, Junhai Xu, Di Jin, Jianhua Tao</i>	
On the Efficacy and Noise-Robustness of Jointly Learned Speech Emotion and Automatic Speech Recognition	1863
<i>Lokesh Bansal, S. Pavankumar Dubagunta, Malolan Chetlur, Pushpak Jagtap, Aravind Ganapathiraju</i>	
Speaking State Decoder with Transition Detection for Next Speaker Prediction.....	1868
<i>Shao-Hao Lu, Yun-Shao Lin, Chi-Chun Lee</i>	
What Are Differences? Comparing DNN and Human by Their Performance and Characteristics in Speaker Age Estimation	1873
<i>Yuki Kitagishi, Naohiro Tawara, Atsunori Ogawa, Ryo Masumura, Taichi Asami</i>	
Effects of Perceived Gender on the Perceived Social Function of Laughter	1878
<i>Joop Arts, Khiet P. Truong</i>	
Implicit Phonetic Information Modeling for Speech Emotion Recognition.....	1883
<i>Tilak Purohit, Bogdan Vlasenko, Mathew Magimai.-Doss</i>	
Computation and Memory Efficient Noise Adaptation of Wav2Vec2.0 for Noisy Speech Emotion Recognition with Skip Connection Adapters.....	1888
<i>Seong-Gyun Leem, Daniel Fulford, Jukka-Pekka Onnela, David Gard, Carlos Busso</i>	
Multi-Level Knowledge Distillation for Speech Emotion Recognition in Noisy Conditions	1893
<i>Yang Liu, Haoqin Sun, Geng Chen, Qingyue Wang, Zhen Zhao, Xugang Lu, Longbiao Wang</i>	
Preference Learning Labels by Anchoring on Consecutive Annotations	1898
<i>Abinay Reddy Naini, Ali N. Salman, Carlos Busso</i>	
Transforming the Embeddings: A Lightweight Technique for Speech Emotion Recognition Tasks	1903
<i>Orchid Chetia Phukan, Arun Balaji Buduru, Rajesh Sharma</i>	
Learning Local to Global Feature Aggregation for Speech Emotion Recognition	1908
<i>Cheng Lu, Hailun Lian, Wenming Zheng, Yuan Zong, Yan Zhao, Sunan Li</i>	
Supervised Contrastive Learning with Nearest Neighbor Search for Speech Emotion Recognition.....	1913
<i>Xuechen Wang, Shiwang Zhao, Yong Qin</i>	

SPEAKER AND LANGUAGE IDENTIFICATION 1

Vietnam-Celeb: A Large-Scale Dataset for Vietnamese Speaker Recognition	1918
<i>Viet Thanh Pham, Xuan Thai Hoa Nguyen, Vu Hoang, Thi Thu Trang Nguyen</i>	
What Can an Accent Identifier Learn? Probing Phonetic and Prosodic Information in a Wav2vec2- Based Accent Identification Model	1923
<i>Mu Yang, Ram C. M. C. Shekar, Okim Kang, John H. L. Hansen</i>	
The 2022 NIST Language Recognition Evaluation.....	1928
<i>Yooyoung Lee, Craig Greenberg, Eliot Godard, Asad A. Butt, Elliot Singer, Trang Nguyen, Lisa Mason, Douglas Reynolds</i>	
Description and Analysis of the KPT System for NIST Language Recognition Evaluation 2022	1933
<i>Salvatore Sarni, Sandro Cumani, Sabato Marco Siniscalchi, Andrea Bottino</i>	
ACA-Net: Towards Lightweight Speaker Verification Using Asymmetric Cross Attention	1938
<i>Jia Qi Yip, Duc-Tuan Truong, Dianwen Ng, Chong Zhang, Yukun Ma, Trung Hieu Nguyen, Chongjia Ni, Shengkui Zhao, Eng Siong Chng, Bin Ma</i>	
Branch-ECAPA-TDNN: A Parallel Branch Architecture to Capture Local and Global Features for Speaker Verification	1943
<i>Jiadi Yao, Chengdong Liang, Zhendong Peng, Binbin Zhang, Xiao-Lei Zhang</i>	
Speaker Verification Across Ages: Investigating Deep Speaker Embedding Sensitivity to Age Mismatch in Enrollment and Test Speech	1948
<i>Vishwanath Pratap Singh, Md Sahidullah, Tomi Kinnunen</i>	
Wavelet Scattering Transform for Improving Generalization in Low-Resourced Spoken Language Identification	1953
<i>Spandan Dey, Premjeet Singh, Goutam Saha</i>	
A Parameter-Efficient Learning Approach to Arabic Dialect Identification with Pre-Trained General-Purpose Speech Model	1958
<i>Srijith Radhakrishnan, Chao-Han Huck Yang, Sumeer Ahmad Khan, Narsis A. Kiani, David Gomez-Cabrero, Jesper N. Tegner</i>	
HABLA: A Dataset of Latin American Spanish Accents for Voice Anti-Spoofing	1963
<i>Pablo Andrés Tamayo Flórez, Rubén Manrique, Bernardo Pereira Nunes</i>	
Self-Supervised Learning Representation Based Accent Recognition with Persistent Accent Memory	1968
<i>Rui Li, Zhiwei Xie, Haihua Xu, Yizhou Peng, Hexin Liu, Hao Huang, Eng Siong Chng</i>	
Extremely Low Bit Quantization for Mobile Speaker Verification Systems Under 1MB Memory	1973
<i>Bei Liu, Haoyu Wang, Yanmin Qian</i>	
Unsupervised Out-Of-Distribution Dialect Detection with Mahalanobis Distance.....	1978
<i>Sourya Dipta Das, Yash Vadi, Abhishek Unnam, Kuldeep Yadav</i>	
Pyannote.audio 2.1 Speaker Diarization Pipeline: Principle, Benchmark, and Recipe	1983
<i>Hervé Bredin</i>	
Model Compression for DNN-Based Speaker Verification Using Weight Quantization.....	1988
<i>Jingyu Li, Wei Liu, Zhaoyang Zhang, Jiong Wang, Tan Lee</i>	

Multi-Resolution Approach to Identification of Spoken Languages and to Improve Overall Language Diarization System Using Whisper Model	1993
<i>Bhavik Vachhani, Dipesh Singh, Rustom Lawyer</i>	
Improving Generalization Ability of Countermeasures for New Mismatch Scenario by Combining Multiple Advanced Regularization Terms	1998
<i>Chang Zeng, Xin Wang, Xiaoxiao Miao, Erica Cooper, Junichi Yamagishi</i>	
Dynamic Fully-Connected Layer for Large-Scale Speaker Verification	2003
<i>Zhida Song, Liang He, Baowei Zhao, Minqiang Xu, Yu Zheng</i>	

SHOW AND TELL: SPEECH TOOLS, SPEECH ENHANCEMENT, SPEECH SYNTHESIS

DeepFilterNet: Perceptually Motivated Real-Time Speech Enhancement	2008
<i>Hendrik Schröter, Alberto N. Escalante-B., Tobias Rosenkranz, Andreas Maier</i>	
Nkululeko: Machine Learning Experiments on Speaker Characteristics Without Programming.....	2010
<i>Felix Burkhardt, Florian Eyben, Björn W. Schuller</i>	
Sp1NY: A Quick and Flexible Speech Visualisation Tool in Python.....	2012
<i>Sébastien Le Maguer, Mark Anderson, Naomi Harte</i>	
Intonation Control for Neural Text-To-Speech Synthesis with Polynomial Models of F0	2014
<i>Niamh Corkey, Johannah O'Mahony, Simon King</i>	
So-To-Speak: An Exploratory Platform for Investigating the Interplay Between Style and Prosody in TTS.....	2016
<i>Éva Székely, Siyang Wang, Joakim Gustafson</i>	
Comparing /b/ and /d/ with a Single Physical Model of the Human Vocal Tract to Visualize Droplets Produced While Speaking.....	2018
<i>Takayuki Arai, Tsukasa Yoshinaga, Akiyoshi Iida</i>	
Show & Tell: Voice Activity Projection and Turn-Taking	2020
<i>Erik Ekstedt, Gabriel Skantze</i>	
Real Time Detection of Soft Voice for Speech Enhancement.....	2022
<i>Hector A. Cordourier, Georg Stemmer, Sinem Aslan, Tobias Bocklet, Himanshu Bhalla</i>	
Data Augmentation for Diverse Voice Conversion in Noisy Environments	2024
<i>Avani Tanna, Michael Saxon, Amr El Abbadi, William Yang Wang</i>	
Application for Real-Time Audio-Visual Speech Enhancement	2026
<i>Mandar Gogate, Kia Dashtipour, Amir Hussain</i>	

SPEECH SYNTHESIS AND VOICE CONVERSION

Mitigating the Exposure Bias in Sentence-Level Grapheme-To-Phoneme (G2P) Transduction	2028
<i>Eunseop Yoon, Hee Suk Yoon, Dhananjaya Gowda, Soohwan Eom, Daehyeok Kim, John Harvill, Heting Gao, Mark Hasegawa-Johnson, Chanwoo Kim, Chang D. Yoo</i>	
Streaming Parrottron for On-Device Speech-To-Speech Conversion.....	2033
<i>Oleg Rybakov, Fadi Biadisy, Xia Zhang, Liyang Jiang, Phoenix Meadowlark, Shivani Agrawal</i>	

Exploiting Emotion Information in Speaker Embeddings for Expressive Text-To-Speech.....	2038
<i>Zein Shaheen, Tasnima Sadekova, Yulia Matveeva, Alexandra Shirshova, Mikhail Kudinov</i>	
E2E-S2S-VC: End-To-End Sequence-To-Sequence Voice Conversion	2043
<i>Takuma Okamoto, Tomoki Toda, Hisashi Kawai</i>	
DC CoMix TTS: An End-To-End Expressive TTS with Discrete Code Collaborated with Mixer.....	2048
<i>Yerin Choi, Myoung-Wan Koo</i>	
Voice Conversion with Just Nearest Neighbors.....	2053
<i>Matthew Baas, Benjamin Van Niekerk, Herman Kamper</i>	
CFVC: Conditional Filtering for Controllable Voice Conversion	2058
<i>Kou Tanaka, Takuhiro Kaneko, Hirokazu Kameoka, Shogo Seki</i>	
DualVC: Dual-Mode Voice Conversion Using Intra-Model Knowledge Distillation and Hybrid Predictive Coding	2063
<i>Ziqian Ning, Yuepeng Jiang, Pengcheng Zhu, Jixun Yao, Shuai Wang, Lei Xie, Mengxiao Bi</i>	
Attention-Based Interactive Disentangling Network for Instance-Level Emotional Voice Conversion.....	2068
<i>Yun Chen, Lingxiao Yang, Qi Chen, Jian-Huang Lai, Xiaohua Xie</i>	
ALO-VC: Any-To-Any Low-Latency One-Shot Voice Conversion.....	2073
<i>Bohan Wang, Damien Ronssin, Milos Cernak</i>	
Evaluating and Reducing the Distance Between Synthetic and Real Speech Distributions	2078
<i>Christoph Minixhofer, Ondrej Klejch, Peter Bell</i>	
Decoupling Segmental and Prosodic Cues of Non-Native Speech Through Vector Quantization	2083
<i>Waris Quamer, Anurag Das, Ricardo Gutierrez-Osuna</i>	
VC-T: Streaming Voice Conversion Based on Neural Transducer	2088
<i>Hiroki Kanagawa, Takafumi Moriya, Yusuke Ijima</i>	

VOLUME 4

Emo-StarGAN: A Semi-Supervised Any-To-Many Non-Parallel Emotion-Preserving Voice Conversion.....	2093
<i>Suhita Ghosh, Arnab Das, Yamini Sinha, Ingo Siegert, Tim Polzehl, Sebastian Stober</i>	
ControlVC: Zero-Shot Voice Conversion with Time-Varying Controls on Pitch and Speed.....	2098
<i>Meiying Chen, Zhiyao Duan</i>	
Reverberation-Controllable Voice Conversion Using Reverberation Time Estimator	2103
<i>Yeonjong Choi, Chao Xie, Tomoki Toda</i>	
Cross-Utterance Conditioned Coherent Speech Editing.....	2108
<i>Cheng Yu, Yang Li, Weiqin Zu, Fanglei Sun, Zheng Tian, Jun Wang</i>	

SPOKEN LANGUAGE TRANSLATION, INFORMATION RETRIEVAL, SUMMARIZATION, RESOURCES, AND EVALUATION 2

MAVD: The First Open Large-Scale Mandarin Audio-Visual Dataset with Depth Information.....	2113
<i>Jianrong Wang, Yuchen Huo, Li Liu, Tianyi Xu, Qi Li, Sen Li</i>	

CN-Celeb-AV: A Multi-Genre Audio-Visual Dataset for Person Recognition	2118
<i>Lantian Li, Xiaolou Li, Haoyu Jiang, Chen Chen, Ruihai Hou, Dong Wang</i>	
Improving Zero-Shot Cross-Domain Slot Filling Via Transformer-Based Slot Semantics Fusion	2123
<i>Yuhang Li, Xiao Wei, Yuke Si, Longbiao Wang, Xiaobao Wang, Jianwu Dang</i>	
Rethinking Transfer and Auxiliary Learning for Improving Audio Captioning Transformer.....	2128
<i>Wooseok Shin, Hyun Joon Park, Jin Sob Kim, Dongwon Kim, Seungjin Lee, Sung Won Han</i>	
Boosting Punctuation Restoration with Data Generation and Reinforcement Learning	2133
<i>Viet Dac Lai, Abel Salinas, Hao Tan, Trung Bui, Quan Tran, Seunghyun Yoon, Hanieh Deilamsalehy, Franck Dernoncourt, Thien Huu Nguyen</i>	
J-ToneNet: A Transformer-Based Encoding Network for Improving Tone Classification in Continuous Speech Via F0 Sequences.....	2138
<i>Yi-Fen Liu, Xiang-Li Lu</i>	
Towards Cross-Language Prosody Transfer for Dialog.....	2143
<i>Jonathan E. Avila, Nigel G. Ward</i>	
Strategies for Improving Low Resource Speech to Text Translation Relying on Pre-Trained ASR Models.....	2148
<i>Santosh Kesiraju, Marek Sarvaš, Tomáš Pavlíček, Cécile Macaire, Alejandro Ciuba</i>	
ITALIC: An Italian Intent Classification Dataset	2153
<i>Alkis Koudounas, Moreno La Quatra, Lorenzo Vaiani, Luca Colomba, Giuseppe Attanasio, Eliana Pastor, Luca Cagliero, Elena Baralis</i>	
Perceptual and Task-Oriented Assessment of a Semantic Metric for ASR Evaluation.....	2158
<i>Janine Rugayan, Giampiero Salvi, Torbjørn Svendsen</i>	
How ChatGPT is Robust for Spoken Language Understanding?	2163
<i>Guangpeng Li, Lu Chen, Kai Yu</i>	
GigaST: A 10,000-Hour Pseudo Speech Translation Corpus.....	2168
<i>Rong Ye, Chengqi Zhao, Tom Ko, Chutong Meng, Tao Wang, Mingxuan Wang, Jun Cao</i>	
Boosting Chinese ASR Error Correction with Dynamic Error Scaling Mechanism.....	2173
<i>Jiaxin Fan, Yong Zhang, Hanzhang Li, Jianzong Wang, Zhitao Li, Sheng Ouyang, Ning Cheng, Jing Xiao</i>	
Crowdsource-Based Validation of the Audio Cocktail as a Sound Browsing Tool	2178
<i>Per Fallgren, Jens Edlund</i>	
PunCantonese: A Benchmark Corpus for Low-Resource Cantonese Punctuation Restoration from Speech Transcripts.....	2183
<i>Yunxiang Li, Pengfei Liu, Xixin Wu, Helen Meng</i>	
Speech-To-Face Conversion Using Denoising Diffusion Probabilistic Models.....	2188
<i>Shuhei Kato, Taiichi Hashimoto</i>	
Inter-Connection: Effective Connection Between Pre-Trained Encoder and Decoder for Speech Translation.....	2193
<i>Yuta Nishikawa, Satoshi Nakamura</i>	

NOVEL TRANSFORMER MODELS FOR ASR

Conner: Streaming Conformer Without Self-Attention for Interactive Voice Assistants	2198
<i>Martin Radfar, Paulina Lyskawa, Brandon Trujillo, Yi Xie, Kai Zhen, Jahn Heymann, Denis Filimonov, Grant P. Strimel, Nathan Susanj, Athanasios Mouchtaris</i>	
Intra-Ensemble: A New Method for Combining Intermediate Outputs in Transformer-Based Automatic Speech Recognition	2203
<i>Dohee Kim, Jieun Choi, Joon-Hyuk Chang</i>	
A Comparative Study on E-Branchformer Vs Conformer in Speech Recognition, Translation, and Understanding Tasks.....	2208
<i>Yifan Peng, Kwangyoun Kim, Felix Wu, Brian Yan, Siddhant Arora, William Chen, Jiyang Tang, Suwon Shon, Prashant Sridhar, Shinji Watanabe</i>	
HyperConformer: Multi-Head HyperMixer for Efficient Speech Recognition	2213
<i>Florian Mai, Juan Zuluaga-Gomez, Titouan Parcollet, Petr Motlicek</i>	
Memory-Augmented Conformer for Improved End-To-End Long-Form ASR.....	2218
<i>Carlos Carvalho, Alberto Abad</i>	
Towards Effective and Compact Contextual Representation for Conformer Transducer Speech Recognition Systems	2223
<i>Mingyu Cui, Jiawen Kang, Jiajun Deng, Xi Yin, Yutao Xie, Xie Chen, Xunying Liu</i>	

SPEAKER RECOGNITION 1

An Enhanced Res2Net with Local and Global Feature Fusion for Speaker Verification	2228
<i>Yafeng Chen, Siqi Zheng, Hui Wang, Luyao Cheng, Qian Chen, Jiajun Qi</i>	
A Study on Visualization of Voiceprint Feature.....	2233
<i>Jian Zhang, Liang He, Xiaochen Guo, Jing Ma</i>	
VoxTube: A Multilingual Speaker Recognition Dataset.....	2238
<i>Ivan Yakovlev, Anton Okhotnikov, Nikita Torgashov, Rostislav Makarov, Yuri Voevodin, Konstantin Simonchik</i>	
Visualizing Data Augmentation in Deep Speaker Recognition	2243
<i>Pengqi Li, Lantian Li, Askar Hamdulla, Dong Wang</i>	

CROSS-LINGUAL AND MULTILINGUAL ASR

Fast and Efficient Multilingual Self-Supervised Pre-Training for Low-Resource Speech Recognition	2248
<i>Zhilong Zhang, Wei Wang, Yanmin Qian</i>	
UniSplice: Universal Cross-Lingual Data Splicing for Low-Resource ASR.....	2253
<i>Wei Wang, Yanmin Qian</i>	
Allophant: Cross-Lingual Phoneme Recognition with Articulatory Attributes	2258
<i>Kevin Glocker, Aaricia Herygers, Munir Georges</i>	
Phonetic-Assisted Multi-Target Units Modeling for Improving Conformer-Transducer ASR System	2263
<i>Li Li, Dongxing Xu, Haoran Wei, Yanhua Long</i>	

Comparison of Multilingual Self-Supervised and Weakly-Supervised Speech Pre-Training for Adaptation to Unseen Languages	2268
<i>Andrew Rouditchenko, Sameer Khurana, Samuel Thomas, Rogerio Feris, Leonid Karlinsky, Hilde Kuehne, David Harwath, Brian Kingsbury, James Glass</i>	
DistilXLSR: A Light Weight Cross-Lingual Speech Representation Model	2273
<i>Haoyu Wang, Siyuan Wang, Wei-Qiang Zhang, Jinfeng Bai</i>	

VOICE CONVERSION

Emotional Voice Conversion with Semi-Supervised Generative Modeling	2278
<i>Hai Zhu, Huayi Zhan, Hong Cheng, Ying Wu</i>	
Diff-HierVC: Diffusion-Based Hierarchical Voice Conversion with Robust Pitch Generation and Masked Prior for Zero-Shot Speaker Adaptation	2283
<i>Ha-Yeong Choi, Sang-Hoon Lee, Seong-Whan Lee</i>	
S2CD: Self-Heuristic Speaker Content Disentanglement for Any-To-Any Voice Conversion.....	2288
<i>Pengfei Wei, Xiang Yin, Chunfeng Wang, Zhonghao Li, Xinghua Qu, Zhiqiang Xu, Zejun Ma</i>	
Flow-VAE VC: End-To-End Flow Framework with Contrastive Loss for Zero-Shot Voice Conversion.....	2293
<i>Le Xu, Rongxiu Zhong, Ying Liu, Huibao Yang, Shilei Zhang</i>	
Automatic Speech Disentanglement for Voice Conversion Using Rank Module and Speech Augmentation	2298
<i>Zhonghua Liu, Shijun Wang, Ning Chen</i>	
End-To-End Zero-Shot Voice Conversion with Location-Variable Convolutions	2303
<i>Wonjune Kang, Mark Hasegawa-Johnson, Deb Roy</i>	

SPEECH AND LANGUAGE IN HEALTH: FROM REMOTE MONITORING TO MEDICAL CONVERSATIONS 2

Classifying Dementia in the Presence of Depression: A Cross-Corpus Study.....	2308
<i>Franziska Braun, Sebastian P. Bayerl, Paula A. Pérez-Toro, Florian Hönig, Hartmut Leffeld, Thomas Hillemacher, Elmar Nöth, Tobias Bocklet, Korbinian Riedhammer</i>	
Exploiting Cross-Domain and Cross-Lingual Ultrasound Tongue Imaging Features for Elderly and Dysarthric Speech Recognition	2313
<i>Shujie Hu, Xurong Xie, Mengzhe Geng, Mingyu Cui, Jiajun Deng, Guinan Li, Tianzi Wang, Helen Meng, Xunying Liu</i>	
Multi-Class Detection of Pathological Speech with Latent Features: How Does it Perform on Unseen Data?.....	2318
<i>Dominik Wagner, Ilja Baumann, Franziska Braun, Sebastian P. Bayerl, Elmar Nöth, Korbinian Riedhammer, Tobias Bocklet</i>	
Responsiveness, Sensitivity and Clinical Utility of Timing-Related Speech Biomarkers for Remote Monitoring of ALS Disease Progression	2323
<i>Hardik Kothare, Michael Neumann, Jackson Liscombe, Jordan Green, Vikram Ramanarayanan</i>	

Use of Speech Impairment Severity for Dysarthric Speech Recognition	2328
<i>Mengzhe Geng, Zengrui Jin, Tianzi Wang, Shujie Hu, Jiajun Deng, Mingyu Cui, Guinan Li, Jianwei Yu, Xurong Xie, Xunying Liu</i>	
MMLung: Moving Closer to Practical Lung Health Estimation Using Smartphones	2333
<i>Mohammed Mosuily, Lindsay Welch, Jagmohan Chauhan</i>	
Investigating the Utility of Synthetic Data for Doctor-Patient Conversation Summarization	2338
<i>Siyuan Chen, Colin A. Grambow, Mojtaba Kadkhodaie Elyaderani, Alireza Sadeghi, Federico Fancellu, Thomas Schaaf</i>	
Non-Uniform Speaker Disentanglement for Depression Detection from Raw Speech Signals	2343
<i>Jinhan Wang, Vijay Ravi, Abeer Alwan</i>	
PoCaPNet: A Novel Approach for Surgical Phase Recognition Using Speech and X-Ray Images.....	2348
<i>Kubilay Can Demir, Tobias Weise, Matthias May, Axel Schmid, Andreas Maier, Seung Hee Yang</i>	
Combining Multiple Multimodal Speech Features into an Interpretable Index Score for Capturing Disease Progression in Amyotrophic Lateral Sclerosis	2353
<i>Michael Neumann, Hardik Kothare, Vikram Ramanarayanan</i>	
The MASCFLICHT Corpus: Face Mask Type and Coverage Area Recognition from Speech	2358
<i>Adria Mallol-Ragolta, Nils Urbach, Shuo Liu, Anton Batliner, Björn W. Schuller</i>	
Towards Reference Speech Characterization for Health Applications	2363
<i>Catarina Botelho, Alberto Abad, Tanja Schultz, Isabel Trancoso</i>	
Automatic Classification of Hypokinetic and Hyperkinetic Dysarthria Based on GMM-Supervectors	2368
<i>Cristian David Ríos-Urrego, Jan Rusz, Elmar Nöth, Juan Rafael Orozco-Arroyave</i>	
Towards Robust Paralinguistic Assessment for Real-World Mobile Health (mHealth) Monitoring: An Initial Study of Reverberation Effects on Speech	2373
<i>Judith Dineley, Ewan Carr, Faith Matcham, Johnny Downs, Richard J. B. Dobson, Thomas F. Quatieri, Nicholas Cummins</i>	

PATHOLOGICAL SPEECH ANALYSIS 1

Multimodal Assessment of Bulbar Amyotrophic Lateral Sclerosis (ALS) Using a Novel Remote Speech Assessment App	2378
<i>Leif Simmatis, Timothy Pommeé, Yana Yunusova</i>	
On the Use of High Frequency Information for Voice Pathology Classification.....	2383
<i>David Martínez, Dayana Ribas, Eduardo Lleida</i>	
Do Phonatory Features Display Robustness to Characterize Parkinsonian Speech Across Corpora?	2388
<i>Anna Favaro, Tianyu Cao, Thomas Thebaud, Jesus Villalba, Ankur Butala, Najim Dehak, Laureano Moro-Velazquez</i>	
Severity Classification of Parkinson's Disease from Speech Using Single Frequency Filtering-Based Features.....	2393
<i>Sudarsana Reddy Kadiri, Manila Kodali, Paavo Alku</i>	

Comparison of Acoustic Measures of Dysphonia in Parkinson's Disease and Huntington's Disease: Effect of Sex and Speaking Task	2398
<i>Michal Šimek, Tomáš Kouba, Michal Novotný, Tereza Tykalová, Jan Ruz</i>	

Alzheimer Disease Classification Through ASR-Based Transcriptions: Exploring the Impact of Punctuation and Pauses	2403
<i>Lucía Gómez-Zaragozá, Simone Wills, Cristian Tejedor-Garcia, Javier Marín-Morales, Mariano Alcañiz, Helmer Strik</i>	

MULTIMODAL SPEECH EMOTION RECOGNITION

LanSER: Language-Model Supported Speech Emotion Recognition	2408
<i>Taesik Gong, Josh Belanich, Krishna Somandepalli, Arsha Nagrani, Brian Eoff, Brendan Jou</i>	

Fine-Tuned RoBERTa Model with a CNN-LSTM Network for Conversational Emotion Recognition	2413
<i>Jiachen Luo, Huy Phan, Joshua Reiss</i>	

Emotion Label Encoding Using Word Embeddings for Speech Emotion Recognition	2418
<i>Eimear Stanley, Eric Demattos, Anita Klementiev, Piotr Ozimek, Georgia Clarke, Michael Berger, Dimitri Palaz</i>	

Discrimination of the Different Intents Carried by the Same Text Through Integrating Multimodal Information	2423
<i>Zhongjie Li, Gaoyan Zhang, Longbiao Wang, Jianwu Dang</i>	

Meta-Domain Adversarial Contrastive Learning for Alleviating Individual Bias in Self-Sentiment Predictions	2428
<i>Zhi Li, Ryu Takeda, Takahiro Hara</i>	

SWRR: Feature Map Classifier Based on Sliding Window Attention and High-Response Feature Reuse for Multimodal Emotion Recognition	2433
<i>Ziping Zhao, Tian Gao, Haishuai Wang, Björn W. Schuller</i>	

SPEECH CODING AND ENHANCEMENT 2

PCNN: A Lightweight Parallel Conformer Neural Network for Efficient Monaural Speech Enhancement	2438
<i>Xinmeng Xu, Weiping Tu, Yuhong Yang</i>	

Exploring the Interactions Between Target Positive and Negative Information for Acoustic Echo Cancellation	2443
<i>Chang Han, Xinmeng Xu, Weiping Tu, Yuhong Yang, Yajie Liu</i>	

Iterative Autoregression: A Novel Trick to Improve Your Low-Latency Speech Enhancement Model	2448
<i>Pavel Andreev, Nicholas Babaev, Azat Saginbaev, Ivan Shchekotov, Aibek Alanov</i>	

A Multi-Dimensional Deep Structured State Space Approach to Speech Enhancement Using Small- Footprint Models	2453
<i>Pin-Jui Ku, Chao-Han Huck Yang, Sabato Siniscalchi, Chin-Hui Lee</i>	

Domain Adaptation for Speech Enhancement in a Large Domain Gap	2458
<i>Lior Frenkel, Jacob Goldberger, Shlomo E. Chazan</i>	

SCP-GAN: Self-Correcting Discriminator Optimization for Training Consistency Preserving Metric GAN on Speech Enhancement Tasks.....	2463
<i>Vasily Zadorozhnyy, Qiang Ye, Kazuhito Koishida</i>	
A Mask Free Neural Network for Monaural Speech Enhancement.....	2468
<i>Liang Liu, Haixin Guan, Jinlong Ma, Wei Dai, Guangyong Wang, Shaowei Ding</i>	
A Training and Inference Strategy Using Noisy and Enhanced Speech as Target for Speech Enhancement Without Clean Speech.....	2473
<i>Li-Wei Chen, Yao-Fei Cheng, Hung-Shin Lee, Yu Tsao, Hsin-Min Wang</i>	
A Simple RNN Model for Lightweight, Low-Compute and Low-Latency Multichannel Speech Enhancement in the Time Domain	2478
<i>Ashutosh Pandey, Ke Tan, Buye Xu</i>	
High Fidelity Speech Enhancement with Band-Split RNN.....	2483
<i>Jianwei Yu, Hangting Chen, Yi Luo, Rongzhi Gu, Chao Weng</i>	
Focus on the Sound Around You: Monaural Target Speaker Extraction Via Distance and Speaker Information.....	2488
<i>Jiuxin Lin, Peng Wang, Heinrich Dinkel, Jun Chen, Zhiyong Wu, Zhiyong Yan, Yongqing Wang, Junbo Zhang, Yujun Wang</i>	
DFSNet: A Steerable Neural Beamformer Invariant to Microphone Array Configuration for Real-Time, Low-Latency Speech Enhancement	2493
<i>Anton Kovalyov, Kashyap Patel, Issa Panahi</i>	
Speaker-Aware Anti-Spoofing.....	2498
<i>Xuechen Liu, Md Sahidullah, Kong Aik Lee, Tomi Kinnunen</i>	
Impact of Residual Noise and Artifacts in Speech Enhancement Errors on Intelligibility of Human and Machine	2503
<i>Shoko Araki, Ayako Yamamoto, Tsubasa Ochiai, Kenichi Arai, Atsunori Ogawa, Tomohiro Nakatani, Toshio Irino</i>	
EffCRN: An Efficient Convolutional Recurrent Network for High-Performance Speech Enhancement	2508
<i>Marvin Sach, Jan Franzen, Bruno Defraene, Kristoff Fluyt, Maximilian Strake, Wouter Tirry, Tim Fingscheidt</i>	
HAD-ANC: A Hybrid System Comprising an Adaptive Filter and Deep Neural Networks for Active Noise Control	2513
<i>Jungphil Park, Jeong-Hwan Choi, Yungyeo Kim, Joon-Hyuk Chang</i>	
MSAF: A Multiple Self-Attention Field Method for Speech Enhancement.....	2518
<i>Minghang Chu, Jing Wang, Yaoyao Ma, Zhiwei Fan, Mengtao Yang, Chao Xu, Zhi Tao, Di Wu</i>	
Ultra Dual-Path Compression for Joint Echo Cancellation and Noise Suppression.....	2523
<i>Hangting Chen, Jianwei Yu, Yi Luo, Rongzhi Gu, Weihua Li, Zhuocheng Lu, Chao Weng</i>	
ABC-KD: Attention-Based-Compression Knowledge Distillation for Deep Learning-Based Noise Suppression	2528
<i>Yixin Wan, Yuan Zhou, Xiulian Peng, Kai-Wei Chang, Yan Lu</i>	
PLCMOS – a Data-Driven Non-Intrusive Metric for the Evaluation of Packet Loss Concealment Algorithms.....	2533
<i>Lorenz Diener, Marju Purin, Sten Sootla, Ando Saabas, Robert Aichner, Ross Cutler</i>	

PHONETICS, PHONOLOGY, AND PROSODY 1

Effects of Meter, Genre and Experience on Pausing, Lengthening and Prosodic Phrasing in German Poetry Reading	2538
<i>Petra Wagner, Simon Betz</i>	
Comparing First Spectral Moment of Australian English /s/ Between Straight and Gay Voices Using Three Analysis Window Sizes.....	2543
<i>Tünde Szalay, John Holik, Duy Duong Nguyen, James Morandini, Catherine J. Madill</i>	
Universal Automatic Phonetic Transcription into the International Phonetic Alphabet	2548
<i>Chihiro Taguchi, Yusuke Sakai, Parisa Haghani, David Chiang</i>	
Voice Twins: Discovering Extremely Similar-Sounding, Unrelated Speakers	2553
<i>Linda Gerlach, Kirsty McDougall, Finnian Kelly, Anil Alexander</i>	
Filling the Population Statistics Gap: Swiss German Reference Data on F0 and Speech Tempo for Forensic Contexts	2558
<i>Hannah Hedegard, Andrea Fröhlich, Fabian Tomaschek, Carina Steiner, Adrian Leemann</i>	
Investigating the Syntax-Discourse Interface in the Phonetic Implementation of Discourse Markers.....	2563
<i>Mathilde Hutin, Liesbeth Degand, Marc Allasonnière-Tang</i>	
Evaluation of a Forensic Automatic Speaker Recognition System with Emotional Speech Recordings.....	2568
<i>Robert Essery, Philip Harrison, Vincent Hughes</i>	
An Outlier Analysis of Vowel Formants from a Corpus Phonetics Pipeline	2573
<i>Emily P. Ahn, Gina-Anne Levow, Richard A. Wright, Eleanor Chodroff</i>	
The Hidden Dance of Phonemes and Visage: Unveiling the Enigmatic Link Between Phonemes and Facial Features.....	2578
<i>Liao Qu, Xianwei Zou, Xiang Li, Yandong Wen, Rita Singh, Bhiksha Raj</i>	
Beatboxing Kick Drum Kinematics	2583
<i>Reed Blaylock, Shrikanth Narayanan</i>	
Effects of Hearing Loss and Amplification on Mandarin Consonant Perception	2588
<i>Huali Zhou, Xianming Bei, Nengheng Zheng, Qinglin Meng</i>	
An Acoustic Analysis of Fricative Variation in Three Accents of English	2593
<i>Roland Adams, Calbert Graham</i>	
Acoustic Cues to Stress Perception in Spanish – a Mismatch Negativity Study	2598
<i>Karolina Bros</i>	
Bulgarian Unstressed Vowel Reduction: Received Views Vs Corpus Findings	2603
<i>Mitko Sabev, Bistra Andreeva, Christoph Gabriel, Jonas Gruenke</i>	
An Investigation of Indian Native Language Phonemic Influences on L2 English Pronunciations	2608
<i>Shelly Jain, Priyanshi Pal, Anil Kumar Vuppala, Prasanta Kumar Ghosh, Chiranjeevi Yarra</i>	
Identifying Stable Sections for Formant Frequency Extraction of French Nasal Vowels Based on Difference Thresholds	2613
<i>Hye-Sook Park, Sunhee Kim</i>	

Evaluation of Delexicalization Methods for Research on Emotional Speech	2618
<i>Nicolas Audibert, Francesca Carbone, Maud Champagne-Lavau, Aurélien Said Housseini, Caterina Petrone</i>	

SPOKEN DIALOG SYSTEMS AND CONVERSATIONAL ANALYSIS 2

Relationship Between Auditory and Semantic Entrainment Using Deep Neural Networks (DNN)	2623
<i>Jay Kejriwal, Štefan Benuš</i>	
Unsupervised Auditory and Semantic Entrainment Models with Deep Neural Networks	2628
<i>Jay Kejriwal, Štefan Benuš, Lina M. Rojas-Barahona</i>	
Prosodic Features Improve Sentence Segmentation and Parsing in English	2633
<i>Elizabeth Nielsen, Mark Steedman, Sharon Goldwater</i>	
Estimation of Listening Response Timing by Generative Model and Parameter Control of Response Substantialness Using Dynamic-Prompt-Tune	2638
<i>Toshiki Muromachi, Yoshinobu Kano</i>	
Parameter Selection for Analyzing Conversations with Autism Spectrum Disorder	2643
<i>Tahiya Chowdhury, Veronica Romero, Amanda Stent</i>	
Efficient Multimodal Neural Networks for Trigger-Less Voice Assistants	2648
<i>Sai Srujana Buddi, Utkarsh Oggy Sarawgi, Tashweena Heeramun, Karan Sawney, Ed Yanosik, Saravana Rathinam, Saurabh Adya</i>	
Rapid Lexical Alignment to a Conversational Agent	2653
<i>Rachel Ostrand, Victor S. Ferreira, David Piorkowski</i>	
Multimodal Turn-Taking Model Using Visual Cues for End-Of-Utterance Prediction in Spoken Dialogue Systems	2658
<i>Fuma Kurata, Mao Saeki, Shinya Fujie, Yoichi Matsuyama</i>	
Audio-Visual Praise Estimation for Conversational Video Based on Synchronization-Guided Multimodal Transformer	2663
<i>Nobukatsu Hojo, Saki Mizuno, Satoshi Kobashikawa, Ryo Masumura, Mana Ithori, Hiroshi Sato, Tomohiro Tanaka</i>	
Improving the Response Timing Estimation for Spoken Dialogue Systems by Reducing the Effect of Speech Recognition Delay	2668
<i>Jin Sakuma, Shinya Fujie, Huaibo Zhao, Tetsunori Kobayashi</i>	
Focus-Attention-Enhanced Crossmodal Transformer with Metric Learning for Multimodal Speech Emotion Recognition	2673
<i>Keulbit Kim, Namhyun Cho</i>	
A Multiple-Teacher Pruning Based Self-Distillation (MT-PSD) Approach to Model Compression for Audio-Visual Wake Word Spotting	2678
<i>Haotian Wang, Jun Du, Hengshun Zhou, Chin-Hui Lee, Yuling Ren, Jiangjiang Zhao</i>	
Abusive Speech Detection in Indic Languages Using Acoustic Features.....	2683
<i>Anika A. Spiesberger, Andreas Triantafyllopoulos, Iosif Tsangko, Björn W. Schuller</i>	
Listening to Silences in Contact Center Conversations Using Textual Cues.....	2688
<i>Digvijay Anil Ingle, Ayush Kumar, Jithendra Vepa</i>	

I Learned Error, I Can Fix It! : A Detector-Corrector Structure for ASR Error Calibration	2693
<i>Heui-Yeen Yeen, Min-Ju Kim, Myoung-Wan Koo</i>	
Verbal and Nonverbal Feedback Signals in Response to Increasing Levels of Miscommunication	2698
<i>Maeva Garnier, Eric Le Ferrand, Fabien Ringeval</i>	
Speech-Based Classification of Defensive Communication: A Novel Dataset and Results	2703
<i>Shahin Amiriparian, Lukas Christ, Regina Kushtanova, Maurice Gerczuk, Alexandra Teynor, Björn W. Schuller</i>	
Quantifying the Perceptual Value of Lexical and Non-Lexical Channels in Speech.....	2708
<i>Sarenne Wallbridge, Peter Bell, Catherine Lai</i>	
Relationships Between Gender, Personality Traits and Features of Multi-Modal Data to Responses to Spoken Dialog Systems Breakdown.....	2713
<i>Kazuya Tsubokura, Yurie Iribe, Norihide Kitaoka</i>	
Speaker-Aware Cross-Modal Fusion Architecture for Conversational Emotion Recognition.....	2718
<i>Huan Zhao, Bo Li, Zixing Zhang</i>	

ANALYSIS OF SPEECH AND AUDIO SIGNALS 2

Blind Estimation of Room Impulse Response from Monaural Reverberant Speech with Segmental Generative Neural Network.....	2723
<i>Zhiheng Liao, Feifei Xiong, Juan Luo, Minjie Cai, Eng Siong Chng, Jinwei Feng, Xionghu Zhong</i>	
Emotion-Aware Audio-Driven Face Animation Via Contrastive Feature Disentanglement.....	2728
<i>Xin Ren, Juan Luo, Xionghu Zhong, Minjie Cai</i>	
Anomalous Sound Detection Based on Sound Separation	2733
<i>Kanta Shimonishi, Kota Dohi, Yohei Kawaguchi</i>	
Random Forest Classification of Breathing Phases from Audio Signals Recorded Using Mobile Devices.....	2738
<i>Vitória S. Fahed, Emer P Doheny, Madeleine M Lowery</i>	
GRAVO: Learning to Generate Relevant Audio from Visual Features with Noisy Online Videos	2743
<i>Youngdo Ahn, Chengyi Wang, Yu Wu, Jong Won Shin, Shujie Liu</i>	
Wav2ToBI: A New Approach to Automatic ToBI Transcription	2748
<i>Wanyue Zhai, Mark Hasegawa-Johnson</i>	
Joint-Former: Jointly Regularized and Locally Down-Sampled Conformer for Semi-Supervised Sound Event Detection	2753
<i>Lijian Gao, Qirong Mao, Ming Dong</i>	
Towards Attention-Based Contrastive Learning for Audio Spoof Detection.....	2758
<i>Chirag Goel, Surya Koppiseti, Ben Colman, Ali Shahriyari, Gaurav Bharaj</i>	
Masked Audio Modeling with CLAP and Multi-Objective Learning	2763
<i>Yifei Xin, Xiulian Peng, Yan Lu</i>	
Few-Shot Open-Set Learning for On-Device Customization of KeyWord Spotting Systems.....	2768
<i>Manuele Rusci, Tinne Tuytelaars</i>	

Self-Supervised Dataset Pruning for Efficient Training in Audio Anti-Spoofing.....	2773
<i>Abdul Hameed Azeemi, Ihsan Ayyub Qazi, Agha Ali Raza</i>	
Semantic Segmentation with Bidirectional Language Models Improves Long-Form ASR	2778
<i>W. Ronny Huang, Hao Zhang, Shankar Kumar, Shuo-Yiin Chang, Tara Sainath</i>	
Multi-Microphone Automatic Speech Segmentation in Meetings Based on Circular Harmonics Features	2783
<i>Théo Mariotte, Anthony Larcher, Silvio Montrésor, Jean-Hugh Thomas</i>	

VOLUME 5

Advanced RawNet2 with Attention-Based Channel Masking for Synthetic Speech Detection	2788
<i>Jing Li, Yanhua Long, Yijie Li, Dongxing Xu</i>	
Insights into End-To-End Audio-To-Score Transcription with Real Recordings: A Case Study with Saxophone Works	2793
<i>Juan Carlos Martínez-Sevilla, María Alfaro-Contreras, Jose J. Valero-Mas, Jorge Calvo-Zaragoza</i>	
Whisper-AT: Noise-Robust Automatic Speech Recognizers Are Also Strong General Audio Event Taggers	2798
<i>Yuan Gong, Sameer Khurana, Leonid Karlinsky, James Glass</i>	
Synthetic Voice Spoofing Detection Based on Feature Pyramid Conformer.....	2803
<i>Jingran Gong, Ning Chen</i>	
Learning a Self-Supervised Domain-Invariant Feature Representation for Generalized Audio Deepfake Detection	2808
<i>Yuankun Xie, Haonan Cheng, Yutian Wang, Long Ye</i>	
Application of Knowledge Distillation to Multi-Task Speech Representation Learning	2813
<i>Mine Kerpici, Van Nguyen, Shuhua Zhang, Erik Visser</i>	
DeCoR: Defy Knowledge Forgetting by Predicting Earlier Audio Codes	2818
<i>Xilin Jiang, Yinghao Aaron Li, Nima Mesgarani</i>	
Variational Classifier for Unsupervised Anomalous Sound Detection Under Domain Generalization	2823
<i>Antonio Almudévar, Alfonso Ortega, Luis Vicente, Antonio Miguel, Eduardo Lleida</i>	
FlexiAST: Flexibility is What AST Needs	2828
<i>Jiu Feng, Mehmet Hamza Erol, Joon Son Chung, Arda Senocak</i>	
MCR-Data2vec 2.0: Improving Self-Supervised Speech Pre-Training Via Model-Level Consistency Regularization	2833
<i>Ji Won Yoon, Seok Min Kim, Nam Soo Kim</i>	
Visually-Aware Audio Captioning with Adaptive Audio-Visual Attention	2838
<i>Xubo Liu, Qiushi Huang, Xinhao Mei, Haohe Liu, Qiuqiang Kong, Jianyuan Sun, Shengchen Li, Tom Ko, Yu Zhang, Lilian H. Tang, Mark D. Plumbley, Volkan Kiliç, Wenwu Wang</i>	

SPEECH CODING: PRIVACY

Masking Kernel for Learning Energy-Efficient Representations for Speaker Recognition and Mobile Health.....	2843
<i>Apiwat Ditthapron, Emmanuel O. Agu, Adam C. Lammert</i>	
ESTimate: A Real-Time Speech Transmission Index Estimator with Speech Enhancement Auxiliary Task Using Self-Attention Feature Pyramid Network	2848
<i>Bajian Xiang, Hongkun Liu, Zedong Wu, Su Shen, Xiangdong Zhang</i>	
Efficient Encoder-Decoder and Dual-Path Conformer for Comprehensive Feature Learning in Speech Enhancement.....	2853
<i>Junyu Wang</i>	
Privacy-Preserving Representation Learning for Speech Understanding.....	2858
<i>Minh Tran, Mohammad Soleymani</i>	
Vocoder Drift in X-Vector-based Speaker Anonymization	2863
<i>Michele Panariello, Massimiliano Todisco, Nicholas Evans</i>	
Malafide: A Novel Adversarial Convolutional Noise Attack Against Deepfake and Spoofing Detection Systems	2868
<i>Michele Panariello, Wanying Ge, Hemlata Tak, Massimiliano Todisco, Nicholas Evans</i>	

ANALYSIS OF NEURAL SPEECH REPRESENTATIONS

Speech Self-Supervised Representation Benchmarking: Are We Doing it Right?	2873
<i>Salah Zaiem, Youcef Kemiche, Titouan Parcollet, Slim Essid, Mirco Ravanelli</i>	
An Extension of Disentanglement Metrics and Its Application to Voice.....	2878
<i>Olivier Zhang, Olivier Le Blouch, Nicolas Gengembre, Damien Lolive</i>	
An Information-Theoretic Analysis of Self-Supervised Discrete Representations of Speech	2883
<i>Badr M. Abdullah, Mohammed Maqsood Shaik, Bernd Möbius, Dietrich Klakow</i>	
SpeechGLUE: How Well Can Self-Supervised Speech Models Capture Linguistic Knowledge?	2888
<i>Takanori Ashihara, Takafumi Moriya, Kohei Matsuura, Tomohiro Tanaka, Yusuke Ijima, Taichi Asami, Marc Delcroix, Yukinori Honma</i>	
Comparison of GIF- And SSL-Based Features in Pathological-Voice Detection.....	2893
<i>Akira Sasou, Yang Chen</i>	
What is Learnt by the LEarnable Front-End (LEAF)? Adapting Per-Channel Energy Normalisation (PCEN) to Noisy Conditions	2898
<i>Hanyu Meng, Vidhyasaharan Sethu, Eliathamby Ambikairajah</i>	

END-TO-END ASR

End-To-End Joint Target and Non-Target Speakers ASR.....	2903
<i>Ryo Masumura, Naoki Makishima, Taiga Yamane, Yoshihiko Yamazaki, Saki Mizuno, Mana Ihori, Mihiro Uchida, Keita Suzuki, Hiroshi Sato, Tomohiro Tanaka, Akihiko Takashima, Satoshi Suzuki, Takafumi Moriya, Nobukatsu Hojo, Atsushi Ando</i>	

Improving Frame-Level Classifier for Word Timings with Non-Peaky CTC in End-To-End Automatic Speech Recognition	2908
<i>Xianzhao Chen, Yist Y. Lin, Kang Wang, Yi He, Zejun Ma</i>	
Joint Autoregressive Modeling of End-To-End Multi-Talker Overlapped Speech Recognition and Utterance-Level Timestamp Prediction	2913
<i>Naoki Makishima, Keita Suzuki, Satoshi Suzuki, Atsushi Ando, Ryo Masumura</i>	
Dual-Path Style Learning for End-To-End Noise-Robust Speech Recognition	2918
<i>Yuchen Hu, Nana Hou, Chen Chen, Eng Siong Chng</i>	
Multi-Pass Training and Cross-Information Fusion for Low-Resource End-To-End Accented Speech Recognition	2923
<i>Xuefei Wang, Yanhua Long, Yijie Li, Haoran Wei</i>	
Text-Only Domain Adaptation for End-To-End ASR Using Integrated Text-To-Mel-Spectrogram Generator	2928
<i>Vladimir Bataev, Roman Korostik, Evgeny Shabalin, Vitaly Lavrukhin, Boris Ginsburg</i>	

SPOKEN LANGUAGE UNDERSTANDING, SUMMARIZATION, AND INFORMATION RETRIEVAL

Leveraging Pretrained ASR Encoders for Effective and Efficient End-To-End Speech Intent Classification and Slot Filling	2933
<i>He Huang, Jagadeesh Balam, Boris Ginsburg</i>	
Relation-Based Counterfactual Data Augmentation and Contrastive Learning for Robustifying Natural Language Inference Models	2938
<i>Heerin Yang, Seung-Won Hwang, Jungmin So</i>	
Transfer Learning from Pre-Trained Language Models Improves End-To-End Speech Summarization.....	2943
<i>Kohei Matsuura, Takanori Ashihara, Takafumi Moriya, Tomohiro Tanaka, Takatomo Kano, Atsunori Ogawa, Marc Delcroix</i>	
Audio Retrieval with WavText5K and CLAP Training	2948
<i>Soham Deshmukh, Benjamin Elizalde, Huaming Wang</i>	
Sequence-Level Knowledge Distillation for Class-Incremental End-To-End Spoken Language Understanding	2953
<i>Umberto Cappellazzo, Muqiao Yang, Daniele Falavigna, Alessio Brutti</i>	
Contrastive Disentangled Learning for Memory-Augmented Transformer.....	2958
<i>Jen-Tzung Chien, Shang-En Li</i>	

INVARIANT AND ROBUST PRE-TRAINED ACOUSTIC MODELS

ProsAudit, a Prosodic Benchmark for Self-Supervised Speech Models	2963
<i>Maureen De Seyssel, Marvin Lavechin, Hadrien Titeux, Arthur Thomas, Gwendal Virlet, Andrea Santos Revilla, Guillaume Wisniewski, Bogdan Ludusan, Emmanuel Dupoux</i>	
Self-Supervised Predictive Coding Models Encode Speaker and Phonetic Information in Orthogonal Subspaces	2968
<i>Oli Danyi Liu, Hao Tang, Sharon Goldwater</i>	

Evaluating Context-Invariance in Unsupervised Speech Representations	2973
<i>Mark Hallap, Emmanuel Dupoux, Ewan Dunbar</i>	
CoBERT: Self-Supervised Speech Representation Learning Through Code Representation Learning	2978
<i>Chutong Meng, Junyi Ao, Tom Ko, Mingxuan Wang, Haizhou Li</i>	
Self-Supervised Fine-Tuning for Improved Content Representations by Speaker-Invariant Clustering	2983
<i>Heng-Jui Chang, Alexander H. Liu, James Glass</i>	
Self-Supervised Acoustic Word Embedding Learning Via Correspondence Transformer Encoder	2988
<i>Jingru Lin, Xianghu Yue, Junyi Ao, Haizhou Li</i>	

PATHOLOGICAL SPEECH ANALYSIS 2

A Pipeline to Evaluate the Effects of Noise on Machine Learning Detection of Laryngeal Cancer	2993
<i>Mary Paterson, James Moor, Luisa Cutillo</i>	
ReCLR: Reference-Enhanced Contrastive Learning of Audio Representation for Depression Detection	2998
<i>Pingyue Zhang, Mengyue Wu, Kai Yu</i>	
Automated Multiple Sclerosis Screening Based on Encoded Speech Representations	3003
<i>José Egas-López, Veronika Svindt, Judit Bóna, Ildikó Hoffmann, Gábor Gosztolya</i>	
Cross-Lingual Features for Alzheimer’s Dementia Detection from Speech.....	3008
<i>Thomas Melistas, Lefteris Kapelonis, Nikos Antoniou, Petros Mitseas, Dimitris Sgouropoulos, Theodoros Giannakopoulos, Athanasios Katsamanis, Shrikanth Narayanan</i>	
Careful Whisper - Leveraging Advances in Automatic Speech Recognition for Robust and Interpretable Aphasia Subtype Classification	3013
<i>Mario Zusaq, Laurin Wagner, Theresa Bloder</i>	
Behavioral Analysis of Pathological Speaker Embeddings of Patients During Oncological Treatment of Oral Cancer	3018
<i>Jenthe Thienpondt, Caroline M. Speksnijder, Kris Demuyneck</i>	

SPEECH SYNTHESIS: REPRESENTATION LEARNING

Adversarial Learning of Intermediate Acoustic Feature for End-To-End Lightweight Text-To-Speech	3023
<i>Hyungchan Yoon, Seyun Um, Changhwan Kim, Hong-Goo Kang</i>	
Adapter-Based Extension of Multi-Speaker Text-To-Speech Model for New Speakers	3028
<i>Cheng-Ping Hsieh, Subhankar Ghosh, Boris Ginsburg</i>	
SALTS: Leveraging Self-Supervised Speech Representations for Improved Text-To-Speech Synthesis.....	3033
<i>Ramanan Sivaguru, Vasista Sai Lodagala, S Umesh</i>	
UnitSpeech: Speaker-Adaptive Speech Synthesis with Untranscribed Data	3038
<i>Heeseung Kim, Sungwon Kim, Jiheum Yeom, Sungroh Yoon</i>	

LightVoc: An Upsampling-Free GAN Vocoder Based on Conformer and Inverse Short-Time Fourier Transform.....	3043
<i>Dinh Son Dang, Tung Lam Nguyen, Bao Thang Ta, Tien Thanh Nguyen, Thi Ngoc Anh Nguyen, Dang Linh Le, Nhat Minh Le, Van Hai Do</i>	

ChatGPT-EDSS: Empathetic Dialogue Speech Synthesis Trained from ChatGPT-Derived Context Word Embeddings	3048
<i>Yuki Saito, Shinnosuke Takamichi, Eiji Iimori, Kentaro Tachibana, Hiroshi Saruwatari</i>	

SPEECH PERCEPTION, PRODUCTION, AND ACQUISITION 1

Human Transcription Quality Improvement.....	3053
<i>Jian Gao, Hanbo Sun, Cheng Cao, Zheng Du</i>	

The Effect of Masking Noise on Listeners' Spectral Tilt Preferences	3058
<i>Olympia Simantiraki, Yannis Pantazis, Martin Cooke</i>	

The Effect of Whistled Vowels on Whistled Word Categorization for Naive Listeners	3063
<i>Anais Tran Ngoc, Fanny Meunier, Julien Meyer</i>	

Automatic Deep Neural Network-Based Segmental Pronunciation Error Detection of L2 English Speech (L1 Bengali).....	3068
<i>Puja Bharati, Sabyasachi Chandra, Shayamal Kumar Das Mandal</i>	

The Effect of Stress on Mandarin Tonal Perception in Continuous Speech for Spanish-Speaking Learners.....	3073
<i>Lixia Hao, Qi Gong, Jinsong Zhang</i>	

Combining Acoustic and Aerodynamic Data Collection: A Perceptual Evaluation of Acoustic Distortions	3078
<i>Amélie Elmerich, Jiayin Gao, Angélique Amélot, Lise Crevier-Buchman, Shinji Maeda</i>	

Estimating Virtual Targets for Lingual Stop Consonants Using General Tau Theory	3083
<i>Benjamin Elie, Alice Turk</i>	

Using Random Forests to Classify Language as a Function of Syllable Timing in Two Groups: Children with Cochlear Implants and with Normal Hearing	3088
<i>Mark Gibson</i>	

An Improved End-To-End Audio-Visual Speech Recognition Model.....	3093
<i>Sheng Yang, Zheng Gong, Jia Kang</i>	

What Influences the Foreign Accent Strength? Phonological and Grammatical Errors in the Perception of Accentedness	3098
<i>Sarah Wesolek, Piotr Gulowski, Joanna Blaszcak, Marzena Zygis</i>	

Investigating the Perception Production Link Through Perceptual Adaptation and Phonetic Convergence.....	3103
<i>Lena-Marie Huttner, Noël Nguyen, Martin J. Pickering</i>	

Emotion Prompting for Speech Emotion Recognition	3108
<i>Xingfa Zhou, Min Li, Lan Yang, Rui Sun, Xin Wang, Huayi Zhan</i>	

Speech-In-Speech Recognition is Modulated by Familiarity to Dialect.....	3113
<i>Jessica L. L. Chin, Elena Talevska, Mark Antoniou</i>	

BASEN: Time-Domain Brain-Assisted Speech Enhancement Network with Convolutional Cross Attention in Multi-Talker Conditions3117
Jie Zhang, Qingtian Xu, Qiu-Shi Zhu, Zhen-Hua Ling

Are Retroflex-To-Dental Sibilant Substitutions in Polish Children's Speech an Example of a Covert Contrast? a Preliminary Acoustic Study 3122
Zuzanna Miodonska, Claartje Levelt, Natalia Mocko, Michal Krecichwost, Agata Sage, Pawel Badura

SPEAKER AND LANGUAGE IDENTIFICATION 2

Reversible Neural Networks for Memory-Efficient Speaker Verification..... 3127
Bei Liu, Yanmin Qian

ECAPA++: Fine-Grained Deep Embedding Learning for TDNN Based Speaker Verification..... 3132
Bei Liu, Yanmin Qian

TO-Rawnet: Improving RawNet with TCN and Orthogonal Regularization for Fake Audio Detection 3137
Chenglong Wang, Jiangyan Yi, Jianhua Tao, Chu Yuan Zhang, Shuai Zhang, Ruibo Fu, Xun Chen

Fooling Speaker Identification Systems with Adversarial Background Music 3142
Chu-Xiao Zuo, Jia-Yi Leng, Wu-Jun Li

Mutual Information-Based Embedding Decoupling for Generalizable Speaker Verification..... 3147
Jianchen Li, Jiqing Han, Shiwen Deng, Tieran Zheng, Yongjun He, Guibin Zheng

Target Active Speaker Detection with Audio-Visual Cues 3152
Yidi Jiang, Ruijie Tao, Zexu Pan, Haizhou Li

Improving End-To-End Neural Diarization Using Conversational Summary Representations..... 3157
Samuel J. Broughton, Lahiru Samarakoon

Phase Perturbation Improves Channel Robustness for Speech Spoofing Countermeasures..... 3162
Yongyi Zang, You Zhang, Zhiyao Duan

Improving Training Datasets for Resource-Constrained Speaker Recognition Neural Networks..... 3167
Pierre-Michel Bousquet, Mickael Rouvier

Instance-Based Temporal Normalization for Speaker Verification..... 3172
Thanathai Lertpetchpun, Ekapol Chuangsuwanich

On the Robustness of Wav2vec 2.0 Based Speaker Recognition Systems 3177
Sergey Novoselov, Galina Lavrentyeva, Anastasia Avdeeva, Vladimir Volokhov, Nikita Khmelev, Artem Akulov, Polina Leonteva

P-Vectors: A Parallel-Coupled TDNN/Transformer Network for Speaker Verification 3182
Xiyuan Wang, Fangyuan Wang, Bo Xu, Liang Xu, Jing Xiao

Group GMM-ResNet for Detection of Synthetic Speech Attacks 3187
Zhenchun Lei, Yan Wen, Yingen Yang, Changhong Liu, Minglei Ma

Robust Training for Speaker Verification Against Noisy Labels..... 3192
Zhihua Fang, Liang He, Hanhan Ma, Xiaochen Guo, Lin Li

Self-Distillation into Self-Attention Heads for Improving Transformer-Based End-To-End Neural Speaker Diarization	3197
<i>Ye-Rin Jeoung, Jeong-Hwan Choi, Ju-Seok Seong, Jehyun Kyung, Joon-Hyuk Chang</i>	
Build a SRE Challenge System: Lessons from VoxSRC 2022 and CNSRC 2022	3202
<i>Zhengyang Chen, Bing Han, Xu Xiang, Houjun Huang, Bei Liu, Yanmin Qian</i>	
Describing the Phonetics in the Underlying Speech Attributes for Deep and Interpretable Speaker Recognition	3207
<i>Imen Ben-Amor, Jean-François Bonastre, Benjamin O'Brien, Pierre-Michel Bousquet</i>	
Range-Based Equal Error Rate for Spoof Localization.....	3212
<i>Lin Zhang, Xin Wang, Erica Cooper, Nicholas Evans, Junichi Yamagishi</i>	
Exploring the English Accent-Independent Features for Speech Emotion Recognition Using Filter and Wrapper-Based Methods for Feature Selection	3217
<i>Nowshin Tabassum, Tasfia Tabassum, Fardin Saad, Tahiya Sultana Safa, Hasan Mahmud, Md. Kamrul Hasan</i>	
Powerset Multi-Class Cross Entropy Loss for Neural Speaker Diarization	3222
<i>Alexis Plaquet, Hervé Bredin</i>	
A Method of Audio-Visual Person Verification by Mining Connections Between Time Series.....	3227
<i>Peiwen Sun, Shanshan Zhang, Zishan Liu, Yougen Yuan, Taotao Zhang, Honggang Zhang, Pengfei Hu</i>	

SPEECH RECOGNITION: ARCHITECTURE, SEARCH, AND LINGUISTIC COMPONENTS

3

A Model for Every User and Budget: Label-Free and Personalized Mixed-Precision Quantization.....	3232
<i>Edward Fish, Umberto Michieli, Mete Ozay</i>	
Modeling Dependent Structure for Utterances in ASR Evaluation	3237
<i>Zhe Liu, Fuchun Peng</i>	
ASR for Low Resource and Multilingual Noisy Code-Mixed Speech.....	3242
<i>Tushar Verma, Atul Shree, Ashutosh Modi</i>	
Accurate and Reliable Confidence Estimation Based on Non-Autoregressive End-To-End Speech Recognition System.....	3247
<i>Xian Shi, Haoneng Luo, Zhifu Gao, Shiliang Zhang, Zhijie Yan</i>	
Combining Multilingual Resources and Models to Develop State-Of-The-Art E2E ASR for Swedish	3252
<i>Lukas Mateju, Jan Nouza, Petr Cerva, Jindrich Zdansky, Frantisek Kynych</i>	
Two Stage Contextual Word Filtering for Context Bias in Unified Streaming and Non-Streaming Transducer	3257
<i>Zhanheng Yang, Sining Sun, Xiong Wang, Yike Zhang, Long Ma, Lei Xie</i>	
Towards Continually Learning New Languages	3262
<i>Ngoc-Quan Pham, Jan Niehues, Alex Waibel</i>	
N-Best T5: Robust ASR Error Correction Using Multiple Input Hypotheses and Constrained Decoding Space.....	3267
<i>Rao Ma, Mark J. F. Gales, Kate M. Knill, Mengjie Qian</i>	

SememeASR: Boosting Performance of End-To-End Speech Recognition Against Domain and Long-Tailed Data Shift with Sememe Semantic Knowledge	3272
<i>Jiaxu Zhu, Changhe Song, Zhiyong Wu, Helen Meng</i>	
MiniStreamer: Enhancing Small Conformer with Chunked-Context Masking for Streaming ASR Applications on the Edge.....	3277
<i>Haris Gulzar, Monikka Roslianna Busto, Takeharu Eda, Katsutoshi Itoyama, Kazuhiro Nakadai</i>	
CoMFLP: Correlation Measure Based Fast Search on ASR Layer Pruning.....	3282
<i>Wei Liu, Zhiyuan Peng, Tan Lee</i>	
Exploration on HuBERT with Multiple Resolution.....	3287
<i>Jiatong Shi, Yun Tang, Hirofumi Inaguma, Hongyu Gong, Juan Pino, Shinji Watanabe</i>	
Quantization-Aware and Tensor-Compressed Training of Transformers for Natural Language Understanding	3292
<i>Zi Yang, Samridhi Choudhary, Siegfried Kunzmann, Zheng Zhang</i>	
Word-Level Confidence Estimation for CTC Models	3297
<i>Burin Naowarat, Thananchai Kongthaworn, Ekapol Chuangsuwanich</i>	
Multilingual Contextual Adapters to Improve Custom Word Recognition in Low-Resource Languages.....	3302
<i>Devang Kulshreshtha, Saket Dingliwal, Brady Houston, Sravan Bodapati</i>	
Unsupervised Active Learning: Optimizing Labeling Cost-Effectiveness for Automatic Speech Recognition	3307
<i>Zhisheng Zheng, Ziyang Ma, Yu Wang, Xie Chen</i>	
4D ASR: Joint Modeling of CTC, Attention, Transducer, and Mask-Predict Decoders.....	3312
<i>Yui Sudo, Shakeel Muhammad, Brian Yan, Jiatong Shi, Shinji Watanabe</i>	
Neural Model Reprogramming with Similarity Based Mapping for Low-Resource Spoken Command Recognition.....	3317
<i>Hao Yen, Pin-Jui Ku, Chao-Han Huck Yang, Hu Hu, Sabato Marco Siniscalchi, Pin-Yu Chen, Yu Tsao</i>	
Language-Specific Boundary Learning for Improving Mandarin-English Code-Switching Speech Recognition	3322
<i>Zhiyun Fan, Linhao Dong, Chen Shen, Zhenlin Liang, Jun Zhang, Lu Lu, Zejun Ma</i>	
Mixture-Of-Expert Conformer for Streaming Multilingual ASR.....	3327
<i>Ke Hu, Bo Li, Tara Sainath, Yu Zhang, Françoise Beaufays</i>	
Lossless 4-Bit Quantization of Architecture Compressed Conformer ASR Systems on the 300-Hr Switchboard Corpus	3332
<i>Zhaoqing Li, Tianzi Wang, Jiajun Deng, Junhao Xu, Shoukang Hu, Xunying Liu</i>	
Compressed MoE ASR Model Based on Knowledge Distillation and Quantization	3337
<i>Yuping Yuan, Zhao You, Shulin Feng, Dan Su, Yanchun Liang, Xiaohu Shi, Dong Yu</i>	

ACOUSTIC MODEL ADAPTATION FOR ASR

- Factorised Speaker-Environment Adaptive Training of Conformer Speech Recognition Systems 3342
Jiajun Deng, Guinan Li, Xurong Xie, Zengrui Jin, Mingyu Cui, Tianzi Wang, Shujie Hu, Mengzhe Geng, Xunying Liu
- Text Only Domain Adaptation with Phoneme Guided Data Splicing for End-To-End Speech Recognition 3347
Wei Wang, Xun Gong, Hang Shao, Dongning Yang, Yanmin Qian
- Cross-Lingual Cross-Age Adaptation for Low-Resource Elderly Speech Emotion Recognition..... 3352
Samuel Cahyawijaya, Holy Lovenia, Willy Chung, Rita Frieske, Zihan Liu, Pascale Fung
- Modular Domain Adaptation for Conformer-Based Streaming ASR 3357
Qijia Li, Bo Li, Dongseong Hwang, Tara Sainath, Pedro M. Mengibar
- Don't Stop Self-Supervision: Accent Adaptation of Speech Representations Via Residual Adapters..... 3362
Anshu Bhatia, Sanchit Sinha, Saket Dingliwal, Karthik Gopalakrishnan, Sravan Bodapati, Katrin Kirchhoff
- SGEM: Test-Time Adaptation for Automatic Speech Recognition Via Sequential-Level Generalized Entropy Minimization 3367
Changhun Kim, Joonhyung Park, Hajin Shim, Eunho Yang

SPEECH SYNTHESIS: EXPRESSIVITY

- A Generative Framework for Conversational Laughter: Its 'Language Model' and Laughter Sound Synthesis..... 3372
Hiroki Mori, Shunya Kimura
- Towards Spontaneous Style Modeling with Semi-Supervised Pre-Training for Conversational Text-To-Speech Synthesis..... 3377
Wei Qin Li, Shun Lei, Qiaochu Huang, Yixuan Zhou, Zhiyong Wu, Shiyin Kang, Helen Meng
- Beyond Style: Synthesizing Speech with Pragmatic Functions..... 3382
Harm Lameris, Joakim Gustafson, Éva Székely
- ECat: An End-To-End Model for Multi-Speaker TTS & Many-To-Many Fine-Grained Prosody Transfer 3387
Ammar Abbas, Sri Karlapati, Bastian Schnell, Penny Karanasou, Marcel Granero Moya, Amith Nagaraj, Ayman Boustati, Nicole Peinelt, Alexis Moinet, Thomas Drugman

MULTI-MODAL SYSTEMS

- BeAts: Bengali Speech Acts Recognition Using Multimodal Attention Fusion..... 3392
Ahana Deb, Sayan Nag, Ayan Mahapatra, Soumitri Chattopadhyay, Aritra Marik, Pijush Kanti Gayen, Shankha Sanyal, Archi Banerjee, Samir Karmakar
- Improving the Gap in Visual Speech Recognition Between Normal and Silent Speech Based on Metric Learning 3397
Sara Kashiwagi, Keitaro Tanaka, Qi Feng, Shigeo Morishima
- Whistle-To-Text: Automatic Recognition of the Silbo Gomero Whistled Language 3402
Agata Jakubiak

A Novel Interpretable and Generalizable Re-Synchronization Model for Cued Speech Based on a Multi-Cuer Corpus	3407
<i>Lufei Gao, Shan Huang, Li Liu</i>	
Visually Grounded Few-Shot Word Acquisition with Fewer Shots.....	3412
<i>Leanne Nortje, Benjamin Van Niekerk, Herman Kamper</i>	
JAMFN: Joint Attention Multi-Scale Fusion Network for Depression Detection	3417
<i>Li Zhou, Zhenyu Liu, Zixuan Shanguan, Xiaoyan Yuan, Yutong Li, Bin Hu</i>	

QUESTION ANSWERING FROM SPEECH

Prompt Guided Copy Mechanism for Conversational Question Answering	3422
<i>Yong Zhang, Zhitao Li, Jianzong Wang, Yiming Gao, Ning Cheng, Fengying Yu, Jing Xiao</i>	
Composing Spoken Hints for Follow-On Question Suggestion in Voice Assistants	3427
<i>Pedro Faustini, Besnik Fetahu, Giuseppe Castellucci, Anjie Fang, Oleg Rokhlenko, Shervin Malmasi</i>	
On Monotonic Aggregation for Open-Domain QA	3432
<i>Sang-Eun Han, Yeonseok Jeong, Seung-Won Hwang, Kyungjae Lee</i>	
Question-Context Alignment and Answer-Context Dependencies for Effective Answer Sentence Selection	3437
<i>Minh Van Nguyen, Kishan KC, Toan Nguyen, Thien Huu Nguyen, Ankit Chadha, Thuy Vu</i>	
Multi-Scale Attention for Audio Question Answering	3442
<i>Guangyao Li, Yixin Xu, Di Hu</i>	
Enhancing Visual Question Answering Via Deconstructing Questions and Explicating Answers	3447
<i>Feilong Chen, Minglun Han, Jing Shi, Shuang Xu, Bo Xu</i>	

MULTI-TALKER METHODS IN SPEECH PROCESSING

SEF-Net: Speaker Embedding Free Target Speaker Extraction Network.....	3452
<i>Bang Zeng, Suo Hongbin, Yulong Wan, Ming Li</i>	
Cascaded Encoders for Fine-Tuning ASR Models on Overlapped Speech	3457
<i>Richard Rose, Oscar Chang, Olivier Siohan</i>	
TokenSplit: Using Discrete Speech Representations for Direct, Refined, and Transcript-Conditioned Speech Separation and Recognition.....	3462
<i>Hakan Erdogan, Scott Wisdom, Xuankai Chang, Zalán Borsos, Marco Tagliasacchi, Neil Zeghidour, John R. Hershey</i>	
Unified Modeling of Multi-Talker Overlapped Speech Recognition and Diarization with a Sidecar Separator.....	3467
<i>Lingwei Meng, Jiawen Kang, Mingyu Cui, Haibin Wu, Xixin Wu, Helen Meng</i>	
Time-Domain Transformer-Based Audiovisual Speaker Separation	3472
<i>Vahid Ahmadi Kalkhorani, Anurag Kumar, Ke Tan, Buye Xu, Deliang Wang</i>	

Multi-Stream Extension of Variational Bayesian HMM Clustering (MS-VBx) for Combined End-To-End and Vector Clustering-Based Diarization.....	3477
<i>Marc Delcroix, Naohiro Tawara, Mireia Diez, Federico Landini, Anna Silnova, Atsunori Ogawa, Tomohiro Nakatani, Lukáš Burget, Shoko Araki</i>	

VOLUME 6

Unsupervised Adaptation with Quality-Aware Masking to Improve Target-Speaker Voice Activity Detection for Speaker Diarization	3482
<i>Shutong Niu, Jun Du, Maokui He, Chin-Hui Lee, Baoxiang Li, Jiakui Li</i>	
BA-SOT: Boundary-Aware Serialized Output Training for Multi-Talker ASR.....	3487
<i>Yuhao Liang, Fan Yu, Yangze Li, Pengcheng Guo, Shiliang Zhang, Qian Chen, Lei Xie</i>	
Improving Label Assignments Learning by Dynamic Sample Dropout Combined with Layer-Wise Optimization in Speech Separation	3492
<i>Chenyang Gao, Yue Gu, Ivan Marsic</i>	
Joint Compensation of Multi-Talker Noise and Reverberation for Speech Enhancement with Cochlear Implants Using One Or More Microphones.....	3497
<i>Clément Gaultier, Tobias Goehring</i>	
Speaker Diarization for ASR Output with T-Vectors: A Sequence Classification Approach.....	3502
<i>Midia Yousefi, Naoyuki Kanda, Dongmei Wang, Zhuo Chen, Xiaofei Wang, Takuya Yoshioka</i>	
GPU-Accelerated Guided Source Separation for Meeting Transcription	3507
<i>Desh Raj, Daniel Povey, Sanjeev Khudanpur</i>	
Overlap Aware Continuous Speech Separation Without Permutation Invariant Training	3512
<i>Linfeng Yu, Wangyou Zhang, Chenda Li, Yanmin Qian</i>	
Weakly-Supervised Speech Pre-Training: A Case Study on Target Speech Recognition.....	3517
<i>Wangyou Zhang, Yanmin Qian</i>	
Directional Speech Recognition for Speaker Disambiguation and Cross-Talk Suppression.....	3522
<i>Ju Lin, Niko Moritz, Ruiming Xie, Kaustubh Kalgaonkar, Christian Fuegen, Frank Seide</i>	
Mixture Encoder for Joint Speech Separation and Recognition.....	3527
<i>Simon Berger, Peter Vieting, Christoph Boeddeker, Ralf Schlüter, Reinhold Haeb-Umbach</i>	

SOCIOPHONETICS

Aberystwyth English Pre-Aspiration in Apparent Time	3532
<i>Miša Michaela Hejná, Adèle Jatteau</i>	
Speech Entrainment in Chinese Story-Style Talk Shows: The Interaction Between Gender and Role	3537
<i>Yanting Sun, Hongwei Ding</i>	
Sociodemographic and Attitudinal Effects on Dialect Speakers' Articulation of the Standard Language: Evidence from German-Speaking Switzerland.....	3542
<i>Carina Steiner, Dieter Studer-Joho, Corinne Lanthemann, Andrin Büchler, Adrian Leemann</i>	
Vowel Normalisation in Latent Space for Sociolinguistics	3547
<i>James Burridge</i>	

SPEAKER AND LANGUAGE DIARIZATION

Attention-Based Encoder-Decoder Network for End-To-End Neural Speaker Diarization with Target Speaker Attractor	3552
<i>Zhengyang Chen, Bing Han, Shuai Wang, Yanmin Qian</i>	
Robust Self Supervised Speech Embeddings for Child-Adult Classification in Interactions Involving Children with Autism	3557
<i>Rimta Lahiri, Tiantian Feng, Rajat Hebbar, Catherine Lord, So Hyun Kim, Shrikanth Narayanan</i>	
The DISPLACE Challenge 2023 - Diarization of SPeaker and LAnguage in Conversational Environments.....	3562
<i>Shikha Baghel, Shreyas Ramoji, Sidharth, Ranjana H, Prachi Singh, Somil Jain, Pratik Roy Chowdhuri, Kaustubh Kulkarni, Swapnil Padhi, Deepu Vijayasenana, Sriram Ganapathy</i>	
Lexical Speaker Error Correction: Leveraging Language Models for Speaker Diarization Error Correction.....	3567
<i>Rohit Paturi, Sundararajan Srinivasan, Xiang Li</i>	
The SpeeD--ZevoTech Submission at DISPLACE 2023	3572
<i>Gabriel Pirlogeanu, Dan Oneata, Alexandru-Lucian Georgescu, Horia Cucu</i>	
End-To-End Neural Speaker Diarization with Absolute Speaker Loss	3577
<i>Chao Wang, Jie Li, Xiang Fang, Jian Kang, Yongxiang Li</i>	

SPEECH EMOTION RECOGNITION 2

A Context-Constrained Sentence Modeling for Deception Detection in Real Interrogation	3582
<i>Ya-Tse Wu, Yuan-Ting Chang, Shao-Hao Lu, Jing-Yi Chuang, Chi-Chun Lee</i>	
MetricAug: A Distortion Metric-Lead Augmentation Strategy for Training Noise-Robust Speech Emotion Recognizer	3587
<i>Ya-Tse Wu, Chi-Chun Lee</i>	
The Co-Use of Laughter and Head Gestures Across Speech Styles	3592
<i>Bogdan Ludusan, Marin Schröer, Martina Rossi, Petra Wagner</i>	
EmotionNAS: Two-Stream Neural Architecture Search for Speech Emotion Recognition	3597
<i>Haiyang Sun, Zheng Lian, Bin Liu, Ying Li, Jianhua Tao, Licai Sun, Cong Cai, Meng Wang, Yuan Cheng</i>	
Pre-Finetuning for Few-Shot Emotional Speech Recognition.....	3602
<i>Maximillian Chen, Zhou Yu</i>	
Integrating Emotion Recognition with Speech Recognition and Speaker Diarisation for Conversations	3607
<i>Wen Wu, Chao Zhang, Philip C. Woodland</i>	
Utility-Preserving Privacy-Enabled Speech Embeddings for Emotion Detection.....	3612
<i>Chandrashekhara Lavania, Sanjiv Das, Xin Huang, Kyu J. Han</i>	
Node-Weighted Graph Convolutional Network for Depression Detection in Transcribed Clinical Interviews	3617
<i>Sergio Burdisso, Esau Villatoro-Tello, Srikanth Madikeri, Petr Motlicek</i>	

Laughter in Task-Based Settings: Whom We Talk to Affects How, When, and How Often We Laugh	3622
<i>Catarina Branco, Isabel Trancoso, Paulo Infante, Khiet P. Truong</i>	
Exploring Downstream Transfer of Self-Supervised Features for Speech Emotion Recognition	3627
<i>Yuanbo Fang, Xiaofen Xing, Xiangmin Xu, Weibin Zhang</i>	
Leveraging Semantic Information for Efficient Self-Supervised Emotion Recognition with Audio- Textual Distilled Models	3632
<i>Danilo De Oliveira, Navin Raj Prabhu, Timo Gerkmann</i>	
Two-Stage Finetuning of Wav2vec 2.0 for Speech Emotion Recognition with ASR and Gender Pretraining	3637
<i>Yuan Gao, Chenhui Chu, Tatsuya Kawahara</i>	
Investigating Acoustic Cues for Multilingual Abuse Detection	3642
<i>Yash Thakran, Vinayak Abrol</i>	
A Novel Frequency Warping Scale for Speech Emotion Recognition	3647
<i>Premjeet Singh, Goutam Saha</i>	
Multi-Scale Temporal Transformer for Speech Emotion Recognition	3652
<i>Zhipeng Li, Xiaofen Xing, Yuanbo Fang, Weibin Zhang, Hengsheng Fan, Xiangmin Xu</i>	
Distant Speech Emotion Recognition in an Indoor Human-Robot Interaction Scenario.....	3657
<i>Nicolás Grágeda, Eduardo Alvarado, Rodrigo Mahu, Carlos Busso, Néstor Becerra Yoma</i>	
A Study on Prosodic Entrainment in Relation to Therapist Empathy in Counseling Conversation	3662
<i>Dehua Tao, Tan Lee, Harold Chui, Sarah Luk</i>	

SHOW AND TELL: LANGUAGE LEARNING AND EDUCATIONAL RESOURCES

A Unified Framework to Improve Learners' Skills of Perception and Production Based on Speech Shadowing and Overlapping	3667
<i>Nobuaki Minematsu, Noriko Nakanishi, Yingxiang Gao, Haitong Sun</i>	
Speak & Improve: L2 English Speaking Practice Tool	3669
<i>Diane Nicholls, Kate M. Knill, Mark J. F. Gales, Anton Ragni, Paul Ricketts</i>	
Measuring Prosody in Child Speech Using SoapBox Fluency API.....	3671
<i>Mauro Nicolao, Brenda McGuirk, Declan Moore, Niall Mullally, Lora Lynn O'Mahony, Emma O'Neill, Amelia C. Kelly</i>	
Teaching Non-Native Sound Contrasts Using Visual Biofeedback.....	3673
<i>Shawn Nissen</i>	
Large-Scale Automatic Audiobook Creation.....	3675
<i>Brendan Walsh, Mark Hamilton, Greg Newby, Xi Wang, Serena Ruan, Sheng Zhao, Lei He, Shaofei Zhang, Eric Dettinger, William T. Freeman, Markus Weimer</i>	
QVoice: Arabic Speech Pronunciation Learning Application.....	3677
<i>Yassine El Kheir, Fouad Khnaisser, Shammur Absar Chowdhury, Hamdy Mubarak, Shazia Afzal, Ahmed M. Ali</i>	
Asking Questions: An Innovative Way to Interact with Oral History Archives	3679
<i>Jan Švec, Martin Bulín, Adam Frémund, Filip Polák</i>	

DisfluencyFixer: A Tool to Enhance Language Learning Through Speech to Speech Disfluency Correction..... 3681
Vineet Bhat, Preethi Jyothi, Pushpak Bhattacharyya

Technology Pipeline for Large Scale Cross-Lingual Dubbing of Lecture Videos into Multiple Indian Languages 3683
Anusha Prakash, Arun Kumar, Ashish Seth, Bhagyashree Mukherjee, Ishika Gupta, Jom Kuriakose, Jordan F, K V Vikram, Mano R Kumar M, Metilda Sagaya Mary, Mohammad Wajahat, Mohana N, Mudit Batra, Navina K, Nihal John George, Nithya Ravi, Pruthwik Mishra, Sudhanshu Srivastava, Vasista Sai Lodagala, Vandan Mujadia, Kada Sai Venkata Vineeth, Vrunda N. Sukhadia, Dipti Sharma, Hema Murthy, Pushpak Bhattacharyya, S Umesh, Rajeev Sangal

MyVoice: Arabic Speech Resource Collaboration Platform..... 3685
Yousseif Elshahawy, Yassine El Kheir, Shammur Absar Chowdhury, Ahmed M. Ali

Personal Primer Prototype 1: Invitation to Make Your Own Embooked Speech-Based Educational Artifact 3687
Daniel D. Hromada, Hyungjoong Kim

ANALYSIS OF SPEECH AND AUDIO SIGNALS 3

Time-Frequency Domain Filter-And-Sum Network for Multi-Channel Speech Separation 3689
Zhewen Deng, Yi Zhou, Hongqing Liu

Audio-Visual Fusion Using Multiscale Temporal Convolutional Attention for Time-Domain Speech Separation..... 3694
Debang Liu, Tianqi Zhang, Mads Græsbøll Christensen, Ying Wei, Zeliang An

An Efficient Speech Separation Network Based on Recurrent Fusion Dilated Convolution and Channel Attention..... 3699
Junyu Wang

Binaural Sound Localization in Noisy Environments Using Frequency-Based Audio Vision Transformer (FAViT)..... 3704
Waradon Phokhinanan, Nicolas Obin, Sylvain Argentieri

Contrastive Learning Based Deep Latent Masking for Music Source Separation..... 3709
Jihyun Kim, Hong-Goo Kang

Speaker Extraction with Detection of Presence and Absence of Target Speakers 3714
Ke Zhang, Marvin Borsdorf, Zexu Pan, Haizhou Li, Yangjie Wei, Yi Wang

PIAVE: A Pose-Invariant Audio-Visual Speaker Extraction Network..... 3719
Qinghua Liu, Meng Ge, Zhizheng Wu, Haizhou Li

Spatial LibriSpeech: An Augmented Dataset for Spatial Audio Learning 3724
Miguel Sarabia, Elena Menyaylenko, Alessandro Toso, Skyler Seto, Zakaria Aldeneh, Shadi Pirhosseinloo, Luca Zappella, Barry-John Theobald, Nicholas Apostoloff, Jonathan Sheaffer

Image-Driven Audio-Visual Universal Source Separation 3729
Chenxing Li, Ye Bai, Yang Wang, Feng Deng, Yuanyuan Zhao, Zhuo Zhang, Xiaorui Wang

Joint Blind Source Separation and Dereverberation for Automatic Speech Recognition Using Delayed-Subsource MNMF with Localization Prior 3734
Mieszko Fras, Marcin Witkowski, Konrad Kowalczyk

SDNet: Stream-Attention and Dual-Feature Learning Network for Ad-Hoc Array Speech Separation.....	3739
<i>Honglong Wang, Chengyun Deng, Yanjie Fu, Meng Ge, Longbiao Wang, Gaoyan Zhang, Jianwu Dang, Fei Wang</i>	
Deeply Supervised Curriculum Learning for Deep Neural Network-Based Sound Source Localization	3744
<i>Min-Sang Baek, Joon-Young Yang, Joon-Hyuk Chang</i>	
Multi-Channel Separation of Dynamic Speech and Sound Events.....	3749
<i>Takuya Fujimura, Robin Scheibler</i>	
Rethinking the Visual Cues in Audio-Visual Speaker Extraction.....	3754
<i>Junjie Li, Meng Ge, Zexu Pan, Rui Cao, Longbiao Wang, Jianwu Dang, Shiliang Zhang</i>	
Using Semi-Supervised Learning for Monaural Time-Domain Speech Separation with a Self-Supervised Learning-Based SI-SNR Estimator.....	3759
<i>Shaoliang Dang, Tetsuya Matsumoto, Yoshinori Takeuchi, Hiroaki Kudo</i>	
Investigation of Training Mute-Expressive End-To-End Speech Separation Networks for an Unknown Number of Speakers.....	3764
<i>Youngwan Kim, Hyungjun Lim, Kiho Yeom, Eunjoon Seo, Hoodong Lee, Stanley Jungkyu Choi, Honglak Lee</i>	
SR-SRP: Super-Resolution Based SRP-PHAT for Sound Source Localization and Tracking	3769
<i>Jae-Heung Cho, Joon-Hyuk Chang</i>	
Dual-Memory Multi-Modal Learning for Continual Spoken Keyword Spotting with Confidence Selection and Diversity Enhancement.....	3774
<i>Zhao Yang, Dianwen Ng, Xizhe Li, Chong Zhang, Rui Jiang, Wei Xi, Yukun Ma, Chongjia Ni, Jizhong Zhao, Bin Ma, Eng Siong Chng</i>	
FN-SSL: Full-Band and Narrow-Band Fusion for Sound Source Localization	3779
<i>Yabo Wang, Bing Yang, Xiaofei Li</i>	
A Neural State-Space Modeling Approach to Efficient Speech Separation	3784
<i>Chen Chen, Chao-Han Huck Yang, Kai Li, Yuchen Hu, Pin-Jui Ku, Eng Siong Chng</i>	
Locate and Beamform: Two-Dimensional Locating All-Neural Beamformer for Multi-Channel Speech Separation	3789
<i>Yanjie Fu, Meng Ge, Honglong Wang, Nan Li, Haoran Yin, Longbiao Wang, Gaoyan Zhang, Jianwu Dang, Chengyun Deng, Fei Wang</i>	
Monaural Speech Separation Method Based on Recurrent Attention with Parallel Branches.....	3794
<i>Xue Yang, Changchun Bao, Xu Zhang, Xianhong Chen</i>	
Ontology-Aware Learning and Evaluation for Audio Tagging.....	3799
<i>Haohe Liu, Qiuqiang Kong, Xubo Liu, Xinhao Mei, Wenwu Wang, Mark D. Plumbley</i>	

SPEECH CODING AND ENHANCEMENT 3

Multi-Dataset Co-Training with Sharpness-Aware Optimization for Audio Anti-Spoofing.....	3804
<i>Hye-Jin Shim, Jee-Weon Jung, Tomi Kinnunen</i>	

Reducing the Prior Mismatch of Stochastic Differential Equations for Diffusion-Based Speech Enhancement	3809
<i>Bunlong Lay, Simon Welker, Julius Richter, Timo Gerkmann</i>	
Complex-Valued Neural Networks for Voice Anti-Spoofing.....	3814
<i>Nicolas M. Müller, Philip Sperl, Konstantin Böttinger</i>	
DeepVQE: Real Time Deep Voice Quality Enhancement for Joint Acoustic Echo Cancellation, Noise Suppression and Dereverberation.....	3819
<i>Nicolae Catalin Ristea, Evgenii Indenbom, Ando Saabas, Tanel Pärnamaa, Jegor Guzhvin, Ross Cutler</i>	
Diffiner: A Versatile Diffusion-Based Generative Refiner for Speech Enhancement.....	3824
<i>Ryosuke Sawata, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Takashi Shibuya, Shusuke Takahashi, Yuki Mitsufuji</i>	
HD-DEMUCS: General Speech Restoration with Heterogeneous Decoders.....	3829
<i>Doyeon Kim, Soo-Whan Chung, Hyewon Han, Youna Ji, Hong-Goo Kang</i>	
MP-SENet: A Speech Enhancement Model with Parallel Denoising of Magnitude and Phase Spectra.....	3834
<i>Ye-Xin Lu, Yang Ai, Zhen-Hua Ling</i>	
TridentSE: Guiding Speech Enhancement with 32 Global Tokens	3839
<i>Dacheng Yin, Zhiyuan Zhao, Chuanxin Tang, Zhiwei Xiong, Chong Luo</i>	
Detection of Cross-Dataset Fake Audio Based on Prosodic and Pronunciation Features	3844
<i>Chenglong Wang, Jiangyan Yi, Jianhua Tao, Chu Yuan Zhang, Shuai Zhang, Xun Chen</i>	
Self-Supervised Learning with Diffusion-Based Multichannel Speech Enhancement for Speaker Verification Under Noisy Conditions	3849
<i>Sandipana Dowerah, Ajinkya Kulkarni, Romain Serizel, Denis Jouviet</i>	
Two-Stage Voice Anonymization for Enhanced Privacy	3854
<i>Francesco Nespola, Daniel Barreda, Jörg Bitzer, Patrick A. Naylor</i>	
Personalized Dereverberation of Speech	3859
<i>Ruilin Xu, Gurunandan Krishnan, Changxi Zheng, Shree K. Nayar</i>	
Weighted Von Mises Distribution-Based Loss Function for Real-Time STFT Phase Reconstruction Using DNN.....	3864
<i>Nguyen Binh Thien, Yukoh Wakabayashi, Yuting Geng, Kenta Iwai, Takano Nishiura</i>	
Deep Multi-Frame Filtering for Hearing Aids.....	3869
<i>Hendrik Schröter, Tobias Rosenkranz, Alberto N. Escalante-B., Andreas Maier</i>	
Aligning Speech Enhancement for Improving Downstream Classification Performance	3874
<i>Yan Xiong, Visar Berisha, Chaitali Chakrabarti</i>	
DNN-Based Parameter Estimation for MVDR Beamforming and Post-Filtering.....	3879
<i>Minseung Kim, Sein Cheong, Jong Won Shin</i>	
FRA-RIR: Fast Random Approximation of the Image-Source Method	3884
<i>Yi Luo, Jianwei Yu</i>	
Rethinking Complex-Valued Deep Neural Networks for Monaural Speech Enhancement.....	3889
<i>Haibin Wu, Ke Tan, Buye Xu, Anurag Kumar, Daniel Wong</i>	

Harmonic Enhancement Using Learnable Comb Filter for Light-Weight Full-Band Speech Enhancement Model	3894
<i>Xiaohuai Le, Tong Lei, Li Chen, Yiqing Guo, Chao He, Cheng Chen, Xianjun Xia, Hua Gao, Yijian Xiao, Piao Ding, Shenyi Song, Jing Lu</i>	

SPOKEN LANGUAGE TRANSLATION, INFORMATION RETRIEVAL, SUMMARIZATION, RESOURCES, AND EVALUATION 3

How Does Pretraining Improve Discourse-Aware Translation?	3899
<i>Zhihong Huang, Longyue Wang, Siyou Liu, Derek F. Wong</i>	
PATCorrect: Non-Autoregressive Phoneme-Augmented Transformer for ASR Error Correction	3904
<i>Ziji Zhang, Zhehui Wang, Rajesh Kamma, Sharanya Eswaran, Narayanan Sadagopan</i>	
Model-Assisted Lexical Tone Evaluation of Three-Year-Old Chinese-Speaking Children by Also Considering Segment Production	3909
<i>Shu-Chuan Tseng, Yi-Fen Liu, Xiang-Li Lu</i>	
Sentence Embedder Guided Utterance Encoder (SEGUE) for Spoken Language Understanding	3914
<i>Yi Xuan Tan, Navonil Majumder, Soujanya Poria</i>	
Joint Time and Frequency Transformer for Chinese Opera Classification	3919
<i>Qiang Li, Beibei Hu</i>	
AdaMS: Deep Metric Learning with Adaptive Margin and Adaptive Scale for Acoustic Word Discrimination	3924
<i>Myunghun Jung, Hoirin Kim</i>	
Investigating Reproducibility at Interspeech Conferences: A Longitudinal and Comparative Perspective	3929
<i>Mohammad Arvan, A. Seza Dogruöz, Natalie Parde</i>	
An Efficient Approach for the Automated Segmentation and Transcription of the People's Speech Corpus	3939
<i>Astik Biswas, Abdelmoumene Boumadane, Stephane Peillon, Gildas Bleas</i>	
Diverse Feature Mapping and Fusion Via Multitask Learning for Multilingual Speech Emotion Recognition	3944
<i>Shi-Wook Lee</i>	
Take the Hint: Improving Arabic Diacritization with Partially-Diacritized Text	3949
<i>Parnia Bahar, Mattia Di Gangi, Nick Rossenbach, Mohammad Zeineldeen</i>	
Low-Resource Cross-Lingual Adaptive Training for Nigerian Pidgin	3954
<i>Pin-Jie Lin, Muhammed Saeed, Ernie Chang, Merel Scholman</i>	
Efficient Adaptation of Spoken Language Understanding Based on End-To-End Automatic Speech Recognition	3959
<i>Eesung Kim, Aditya Jajodia, Cindy Tseng, Divya Neelagiri, Taeyeon Ki, Vijendra Raj Apsingekar</i>	
PhonMatchNet: Phoneme-Guided Zero-Shot Keyword Spotting for User-Defined Keywords	3964
<i>Yong-Hyeok Lee, Namhyun Cho</i>	

Mix Before Align: Towards Zero-Shot Cross-Lingual Sentiment Analysis Via Soft-Mix and Multi-View Learning	3969
<i>Zhihong Zhu, Xuxin Cheng, Dongsheng Chen, Zhiqi Huang, Hongxiang Li, Yuexian Zou</i>	
AlignAtt: Using Attention-Based Audio-Translation Alignments as a Guide for Simultaneous Speech Translation	3974
<i>Sara Papi, Marco Turchi, Matteo Negri</i>	
Incremental Blockwise Beam Search for Simultaneous Speech Translation with Controllable Quality-Latency Tradeoff	3979
<i>Peter Polák, Brian Yan, Shinji Watanabe, Alex Waibel, Ondrej Bojar</i>	
Zambezi Voice: A Multilingual Speech Corpus for Zambian Languages	3984
<i>Claytone Sikasote, Kalinda Siaminwe, Stanly Mwape, Bangiwe Zulu, Mofya Phiri, Martin Phiri, David Zulu, Mayumbo Nyirenda, Antonios Anastasopoulos</i>	

ANTI-SPOOFING FOR SPEAKER VERIFICATION

Towards Single Integrated Spoofing-Aware Speaker Verification Embeddings.....	3989
<i>Sung Hwan Mun, Hye-Jin Shim, Hemlata Tak, Xin Wang, Xuechen Liu, Md Sahidullah, Myeonghun Jeong, Min Hyun Han, Massimiliano Todisco, Kong Aik Lee, Junichi Yamagishi, Nicholas Evans, Tomi Kinnunen, Nam Soo Kim, Jee-Weon Jung</i>	
Pseudo-Siamese Network Based Timbre-Reserved Black-Box Adversarial Attack in Speaker Identification	3994
<i>Qing Wang, Jixun Yao, Ziqian Wang, Pengcheng Guo, Lei Xie</i>	
Betray Oneself: A Novel Audio DeepFake Detection Model Via Mono-To-Stereo Conversion	3999
<i>Rui Liu, Jinhua Zhang, Guanglai Gao, Haizhou Li</i>	
Robust Audio Anti-Spoofing Countermeasure with Joint Training of Front-End and Back-End Models	4004
<i>Xingming Wang, Bang Zeng, Suo Hongbin, Yulong Wan, Ming Li</i>	
Improved DeepFake Detection Using Whisper Features.....	4009
<i>Piotr Kawa, Marcin Plata, Michal Czuba, Piotr Szymanski, Piotr Syga</i>	
DoubleDeceiver: Deceiving the Speaker Verification System Protected by Spoofing Countermeasures	4014
<i>Mengao Zhang, Ke Xu, Hao Li, Lei Wang, Chengfang Fang, Jie Shi</i>	

SPEECH CODING: INTELLIGIBILITY

On Training a Neural Residual Acoustic Echo Suppressor for Improved ASR.....	4019
<i>Sankaran Panchapagesan, Turaj Zakizadeh Shabestary, Arun Narayanan</i>	
Extending DNN-Based Multiplicative Masking to Deep Subband Filtering for Improved Dereverberation	4024
<i>Jean-Marie Lemerrier, Julian Tobergte, Timo Gerkmann</i>	
UnSE: Unsupervised Speech Enhancement Using Optimal Transport.....	4029
<i>Wenbin Jiang, Fei Wen, Yifan Zhang, Kai Yu</i>	

MC-SpEx: Towards Effective Speaker Extraction with Multi-Scale Interfusion and Conditional Speaker Modulation.....	4034
<i>Jun Chen, Wei Rao, Zilin Wang, Jiuxin Lin, Yukai Ju, Shulin He, Yannan Wang, Zhiyong Wu</i>	
Causal Signal-Based DCCRN with Overlapped-Frame Prediction for Online Speech Enhancement	4039
<i>Juliita Bartolewska, Stanislaw Kacprzak, Konrad Kowalczyk</i>	
Gesper: A Restoration-Enhancement Framework for General Speech Reconstruction	4044
<i>Wenzhe Liu, Yupeng Shi, Jun Chen, Wei Rao, Shulin He, Andong Li, Yannan Wang, Zhiyong Wu</i>	

RESOURCES FOR SPOKEN LANGUAGE PROCESSING

Multimodal Personality Traits Assessment (MuPTA) Corpus: The Impact of Spontaneous and Read Speech	4049
<i>Elena Ryumina, Dmitry Ryumin, Maxim Markitantov, Heysem Kaya, Alexey Karpov</i>	
MOCKS 1.0: Multilingual Open Custom Keyword Spotting Testset.....	4054
<i>Mikolaj Pudo, Mateusz Wosik, Adam Cieslak, Justyna Krzywdziak, Bozena Lukasiak, Artur Janicki</i>	
MD3: The Multi-Dialect Dataset of Dialogues	4059
<i>Jacob Eisenstein, Vinodkumar Prabhakaran, Clara Rivera, Dorottya Demszky, Devyani Sharma</i>	
MuAViC: A Multilingual Audio-Visual Corpus for Robust Speech Recognition and Robust Speech-To-Text Translation.....	4064
<i>Mohamed Anwar, Bowen Shi, Vedanuj Goswami, Wei-Ning Hsu, Juan Pino, Changhan Wang</i>	
Thai Dialect Corpus and Transfer-Based Curriculum Learning Investigation for Dialect Automatic Speech Recognition	4069
<i>Artit Suwanbandit, Burin Naowarat, Orathai Sangpetch, Ekapol Chuangsuwanich</i>	
HK-LegiCoST: Leveraging Non-Verbatim Transcripts for Speech Translation.....	4074
<i>Cihan Xiao, Henry Li Xinyuan, Jinyi Yang, Dongji Gao, Matthew Wiesner, Kevin Duh, Sanjeev Khudanpur</i>	

NEW COMPUTATIONAL STRATEGIES FOR ASR TRAINING AND INFERENCE

A Metric-Driven Approach to Conformer Layer Pruning for Efficient ASR Inference.....	4079
<i>Dhanush Bekal, Karthik Gopalakrishnan, Karel Mundnich, Srikanth Ronanki, Sravan Bodapati, Katrin Kirchhoff</i>	
Distillation Strategies for Discriminative Speech Recognition Rescoring	4084
<i>Prashanth Gurunath Shivakumar, Jari Kolehmainen, Yile Gu, Ankur Gandhe, Ariya Rastrow, Ivan Bulyko</i>	
Another Point of View on Visual Speech Recognition.....	4089
<i>Baptiste Pouthier, Laurent Pilati, Giacomo Valenti, Charles Bouveyron, Frederic Precioso</i>	
RASR2: The RWTH ASR Toolkit for Generic Sequence-To-Sequence Speech Recognition.....	4094
<i>Wei Zhou, Eugen Beck, Simon Berger, Ralf Schlüter, Hermann Ney</i>	
Streaming Speech-To-Confusion Network Speech Recognition	4099
<i>Denis Filimonov, Prabhat Pandey, Ariya Rastrow, Ankur Gandhe, Andreas Stolcke</i>	

Accurate and Structured Pruning for Efficient Automatic Speech Recognition.....	4104
<i>Huiqiang Jiang, Li Lyna Zhang, Yuang Li, Yu Wu, Shijie Cao, Ting Cao, Yuqing Yang, Jinyu Li, Mao Yang, Lili Qiu</i>	

MERLION CCS CHALLENGE: MULTILINGUAL EVERYDAY RECORDINGS - LANGUAGE IDENTIFICATION ON CODE-SWITCHED CHILD-DIRECTED SPEECH

MERLion CCS Challenge: A English-Mandarin Code-Switching Child-Directed Speech Corpus for Language Identification and Diarization	4109
<i>Victoria Y. H. Chua, Hexin Liu, Leibny Paola Garcia, Fei Ting Woon, Jinyi Wong, Xiangyu Zhang, Sanjeev Khudanpur, Andy W. H. Khong, Justin Dauwels, Suzy J. Styles</i>	
Spoken Language Identification System for English-Mandarin Code-Switching Child-Directed Speech	4114
<i>Shashi Kant Gupta, Sushant Hiray, Prashant Kukde</i>	
Improving Wav2vec2-Based Spoken Language Identification by Learning Phonological Features.....	4119
<i>Mostafa Shahin, Zheng Nan, Vidhyasaharan Sethu, Beena Ahmed</i>	
Language Identification Networks for Multilingual Everyday Recordings.....	4124
<i>Kiran Praveen, Balaji Radhakrishnan, Kamini Sabu, Abhishek Pandey, Mahaboob Ali Basha Shaik</i>	
Investigating Model Performance in Language Identification: Beyond Simple Error Statistics	4129
<i>Suzy J. Styles, Victoria Y. H. Chua, Fei Ting Woon, Hexin Liu, Leibny Paola Garcia, Sanjeev Khudanpur, Andy W. H. Khong, Justin Dauwels</i>	

HEALTH-RELATED SPEECH ANALYSIS

Classification of Vocal Intensity Category from Speech Using the Wav2vec2 and Whisper Embeddings.....	4134
<i>Manila Kodali, Sudarsana Reddy Kadiri, Paavo Alku</i>	
The Effect of Clinical Intervention on the Speech of Individuals with PTSD: Features and Recognition Performances.....	4139
<i>Alexander Kathan, Andreas Triantafyllopoulos, Shahin Amiriparian, Sabrina Milkus, Alexander Gebhard, Jonas Hohmann, Pauline Muderlak, Jürgen Schottdorf, Björn W. Schuller, Richard Musil</i>	
Analysis and Automatic Prediction of Exertion from Speech: Contrasting Objective and Subjective Measures Collected While Running	4144
<i>Andreas Triantafyllopoulos, Alexander Gebhard, Alexander Kathan, Maurice Gerczuk, Shahin Amiriparian, Björn W. Schuller</i>	
The Androids Corpus: A New Publicly Available Benchmark for Speech Based Depression Detection	4149
<i>Fuxiang Tao, Anna Esposito, Alessandro Vinciarelli</i>	
Comparing Hand-Crafted Features to Spectrograms for Autism Severity Estimation	4154
<i>Marina Eni, Ilan Dinstein, Yaniv Zigel</i>	
Acoustic Characteristics of Depression in Older Adults' Speech: The Role of Covariates	4159
<i>Carmen Mijnders, Esther Janse, Paul Naarding, Khiet P. Truong</i>	

AUTOMATIC AUDIO CLASSIFICATION AND AUDIO CAPTIONING

- Dual Transformer Decoder Based Features Fusion Network for Automated Audio Captioning..... 4164
Jianyuan Sun, Xubo Liu, Xinhao Mei, Volkan Kiliç, Mark D. Plumbley, Wenwu Wang
- Adapting a ConvNeXt Model to Audio Classification on AudioSet 4169
Thomas Pellegrini, Ismail Khalfaoui-Hassani, Etienne Labbé, Timothée Masquelier
- Few-Shot Class-Incremental Audio Classification Using Stochastic Classifier 4174
Yanxiong Li, Wenchang Cao, Jialong Li, Wei Xie, Qianhua He

VOLUME 7

- Enhance Temporal Relations in Audio Captioning with Sound Event Detection..... 4179
Zeyu Xie, Xuenan Xu, Mengyue Wu, Kai Yu

SPEECH PERCEPTION, PRODUCTION, AND ACQUISITION 2

- First Language Effects on Second Language Perception: Evidence from English Low-Vowel Nasal Sequences Perceived by L1 Mandarin Chinese Listeners 4184
Sijia Zhang
- Motor Control Similarity Between Speakers Saying “A Souk” Using Inverse Atlas Tongue Modeling 4189
Ursa Maity, Fangxu Xing, Jerry Prince, Maureen Stone, El Fakhri Georges, Jonghye Woo, Sidney Fels
- Assessing Phrase Break of ESL Speech with Pre-Trained Language Models and Large Language Models..... 4194
Zhiyi Wang, Shaoguang Mao, Wenshan Wu, Yan Xia, Yan Deng, Jonathan Tien
- A Relationship Between Vocal Fold Vibration and Droplet Production 4199
Tsukasa Yoshinaga, Takayuki Arai, Akiyoshi Iida
- Audio, Visual and Audiovisual Intelligibility of Vowels Produced in Noise..... 4204
Maeva Garnier
- Optimal Control of Speech with Context-Dependent Articulatory Targets 4209
Benjamin Elie, Juraj Šimko, Alice Turk
- Computational Modeling of Auditory Brainstem Responses Derived from Modified Speech..... 4214
Tzu-Han Zoe Cheng, Paul Calamia
- Leveraging Label Information for Multimodal Emotion Recognition 4219
Peiying Wang, Sunlu Zeng, Junqing Chen, Lu Fan, Meng Chen, Youzheng Wu, Xiaodong He
- Improving End-To-End Modeling for Mandarin-English Code-Switching Using Lightweight Switch-Routing Mixture-Of-Experts 4224
Fengyun Tan, Chaofeng Feng, Tao Wei, Shuai Gong, Jinqiang Leng, Wei Chu, Jun Ma, Shaojun Wang, Jing Xiao
- Frequency Patterns of Individual Speaker Characteristics at Higher and Lower Spectral Ranges..... 4229
Zhao Zhang, Ju Zhang, Ziyu Zhu, Yujie Chi, Kiyoshi Honda, Jianguo Wei

Adaptation to Predictive Prosodic Cues in Non-Native Standard Dialect.....	4234
<i>Sabine Gosselke Berthelsen</i>	
Head Movements in Two- And Four-Person Interactive Conversational Tasks in Noisy and Moderately Reverberant Conditions.....	4239
<i>Alan Archer-Boyd, Rainer Martin</i>	
Second Language Identification of Vietnamese Tones by Native Mandarin Learners	4244
<i>Juqiang Chen, Ailing Qin, Hui Chang, Hua Chen</i>	
Nasal Vowel Production and Grammatical Processing in French-Speaking Children with Cochlear Implants and Normal-Hearing Peers.	4249
<i>Sophie Fagniard, Véronique Delvaux, Brigitte Charlier, Bernard Harmegnies, Anne Huberlant, Myriam Piccaluga, Kathy Huet</i>	
Emotion Classification with EEG Responses Evoked by Emotional Prosody of Speech	4254
<i>Zechen Zhang, Xihong Wu, Jing Chen</i>	
L2-Mandarin Regional Accent Variability During Mandarin Tone-Word Training Facilitates English Listeners' Subsequent Tone Categorizations	4259
<i>Yanping Li, Michael D. Tyler, Denis Burnham, Catherine T. Best</i>	
HumanDiffusion: Diffusion Model Using Perceptual Gradients.....	4264
<i>Yota Ueda, Shinnosuke Takamichi, Yuki Saito, Norihiro Takamune, Hiroshi Saruwatari</i>	
Queer Events, Relationships, and Sports: Does Topic Influence Speakers' Acoustic Expression of Sexual Orientation?	4269
<i>Sven Kachel, Manuel Pöhlmann, Christine Nussbaum</i>	

SPEECH SYNTHESIS

Epoch-Based Spectrum Estimation for Speech.....	4274
<i>Jón Guðnason, Guolin Fang, Mike Brookes</i>	
OverFlow: Putting Flows on Top of Neural Transducers for Better TTS.....	4279
<i>Shivam Mehta, Ambika Kirkland, Harm Lameris, Jonas Beskow, Éva Székely, Gustav Eje Henter</i>	
ADAPTERMIX: Exploring the Efficacy of Mixture of Adapters for Low-Resource TTS Adaptation.....	4284
<i>Ambuj Mehrish, Abhinav Ramesh Kashyap, Li Yingting, Navonil Majumder, Soujanya Poria</i>	
Prior-Free Guided TTS: An Improved and Efficient Diffusion-Based Text-Guided Speech Synthesis.....	4289
<i>Won-Gook Choi, So-Jeong Kim, Taeho Kim, Joon-Hyuk Chang</i>	
UnDiff: Unsupervised Voice Restoration with Unconditional Diffusion Model.....	4294
<i>Anastasiia Iashchenko, Pavel Andreev, Ivan Shchekotov, Nicholas Babaev, Dmitry Vetrov</i>	
Pruning Self-Attention for Zero-Shot Multi-Speaker Text-To-Speech.....	4299
<i>Hyungchan Yoon, Changhwan Kim, Eunwoo Song, Hyun-Wook Yoon, Hong-Goo Kang</i>	
Interpretable Style Transfer for Text-To-Speech with ControlVAE and Diffusion Bridge.....	4304
<i>Wenhao Guan, Tao Li, Yishuang Li, Hukai Huang, Qingyang Hong, Lin Li</i>	
Towards Robust FastSpeech 2 by Modelling Residual Multimodality.....	4309
<i>Fabian Kögel, Bac Nguyen, Fabien Cardinaux</i>	

Real Time Spectrogram Inversion on Mobile Phone.....	4314
<i>Oleg Rybakov, Marco Tagliasacchi, Yunpeng Li, Liyang Jiang, Xia Zhang, Fadi Biadsy</i>	
Automatic Tuning of Loss Trade-Offs Without Hyper-Parameter Search in End-To-End Zero-Shot Speech Synthesis	4319
<i>Seongyeon Park, Bohyung Kim, Tae-Hyun Oh</i>	
A Low-Resource Pipeline for Text-To-Speech from Found Data with Application to Scottish Gaelic	4324
<i>Dan Wells, Korin Richmond, William Lamb</i>	
Self-Supervised Solution to the Control Problem of Articulatory Synthesis.....	4329
<i>Paul K. Krug, Peter Birkholz, Branislav Gerazov, Daniel R. Van Niekerk, Anqi Xu, Yi Xu</i>	
Hierarchical Timbre-Cadence Speaker Encoder for Zero-Shot Speech Synthesis	4334
<i>Joun Yeop Lee, Jae-Sung Bae, Seongkyu Mun, Jihwan Lee, Ji-Hyun Lee, Hoon-Young Cho, Chanwoo Kim</i>	
ZET-Speech: Zero-Shot Adaptive Emotion-Controllable Text-To-Speech Synthesis with Diffusion and Style-Based Models	4339
<i>Minki Kang, Wooseok Han, Sung Ju Hwang, Eunho Yang</i>	
Improving WaveRNN with Heuristic Dynamic Blending for Fast and High-Quality GPU Vocoding	4344
<i>Muyang Du, Chuan Liu, Jiaying Qi, Junjie Lai</i>	
Intelligible Lip-To-Speech Synthesis with Speech Units	4349
<i>Jeongsoo Choi, Minsu Kim, Yong Man Ro</i>	
Parameter-Efficient Learning for Text-To-Speech Accent Adaptation	4354
<i>Li-Jen Yang, Chao-Han Huck Yang, Jen-Tzung Chien</i>	
Controlling Formant Frequencies with Neural Text-To-Speech for the Manipulation of Perceived Speaker Age.....	4359
<i>Ziya Khan, Lovisa Wihlborg, Cassia Valentini-Botinhao, Oliver Watts</i>	
FastFit: Towards Real-Time Iterative Neural Vocoder by Replacing U-Net Encoder with Multiple STFTs	4364
<i>Won Jang, Dan Lim, Heayoung Park</i>	
ISTFTNet2: Faster and More Lightweight iSTFT-Based Neural Vocoder Using 1D-2D CNN	4369
<i>Takuhiko Kaneko, Hirokazu Kameoka, Kou Tanaka, Shogo Seki</i>	
VITS2: Improving Quality and Efficiency of Single-Stage Text-To-Speech with Adversarial Learning and Architecture Design.....	4374
<i>Jungil Kong, Jihoon Park, Beomjeong Kim, Jeongmin Kim, Dohee Kong, Sangjin Kim</i>	
Controlling Multi-Class Human Vocalization Generation Via a Simple Segment-Based Labeling Scheme	4379
<i>Hieu-Thi Luong, Junichi Yamagishi</i>	
 <u>SPEECH RECOGNITION: SIGNAL PROCESSING, ACOUSTIC MODELING, ROBUSTNESS, ADAPTATION 4</u>	
Vistaar: Diverse Benchmarks and Training Sets for Indian Language ASR.....	4384
<i>Kaushal Bhogale, Sai Sundaresan, Abhigyan Raman, Tahir Javed, Mitesh M. Khapra, Pratyush Kumar</i>	

Domain Adaptive Self-Supervised Training of Automatic Speech Recognition	4389
<i>Cong-Thanh Do, Rama Doddipatla, Mohan Li, Thomas Hain</i>	
There is More than One Kind of Robustness: Fooling Whisper with Adversarial Examples.....	4394
<i>Raphael Olivier, Bhiksha Raj</i>	
MT-SLVR: Multi-Task Self-Supervised Learning for Transformation In(Variant) Representations.....	4399
<i>Calum Heggan, Tim Hospedales, Sam Budgett, Mehrdad Yaghoobi</i>	
Reducing Barriers to Self-Supervised Learning: HuBERT Pre-Training with Academic Compute.....	4404
<i>William Chen, Xuankai Chang, Yifan Peng, Zhaocheng Ni, Soumi Maiti, Shinji Watanabe</i>	
Blank-Regularized CTC for Frame Skipping in Neural Transducer.....	4409
<i>Yifan Yang, Xiaoyu Yang, Liyong Guo, Zengwei Yao, Wei Kang, Fangjun Kuang, Long Lin, Xie Chen, Daniel Povey</i>	
The Tag-Team Approach: Leveraging CLS and Language Tagging for Enhancing Multilingual ASR.....	4414
<i>Kaousheik Jayakumar, Vrunda N. Sukhadia, A Arunkumar, S Umesh</i>	
Improving RNN-Transducers with Acoustic LookAhead	4419
<i>Vinit S. Unni, Ashish Mittal, Preethi Jyothi, Sunita Sarawagi</i>	
Everyone Has an Accent.....	4424
<i>Nina Markl, Catherine Lai</i>	
Some Voices Are Too Common: Building Fair Speech Recognition Systems Using the CommonVoice Dataset	4428
<i>Lucas Maison, Yannick Estève</i>	
Information Magnitude Based Dynamic Sub-Sampling for Speech-To-Text.....	4433
<i>Yuhao Zhang, Chenghao Gao, Kaiqi Kou, Chen Xu, Tong Xiao, Jingbo Zhu</i>	

KEYNOTE 3

What's in a Rise? the Relevance of Intonation for Attention Orienting	4438
<i>Martine Grice</i>	

SPEECH SYNTHESIS: CONTROLLABILITY AND ADAPTATION

HierVST: Hierarchical Adaptive Zero-Shot Voice Style Transfer	4439
<i>Sang-Hoon Lee, Ha-Yeong Choi, Hyung-Seok Oh, Seong-Whan Lee</i>	
VISinger2: High-Fidelity End-To-End Singing Voice Synthesis Enhanced by Digital Signal Processing Synthesizer	4444
<i>Yongmao Zhang, Heyang Xue, Hanzhao Li, Lei Xie, Tingwei Guo, Ruixiong Zhang, Caixia Gong</i>	
EdenTTS: A Simple and Efficient Parallel Text-To-Speech Architecture with Collaborative Duration-Alignment Learning	4449
<i>Youneng Ma, Junyi He, Meimei Wu, Guangyue Hu, Haojun Fei</i>	
Generalizable Zero-Shot Speaker Adaptive Speech Synthesis with Disentangled Representations.....	4454
<i>Wenbin Wang, Yang Song, Sanjay Jha</i>	

Speech Inpainting: Context-Based Speech Synthesis Guided by Video.....	4459
<i>Juan Felipe Montesinos, Daniel Michelsanti, Gloria Haro, Zheng-Hua Tan, Jesper Jensen</i>	
STEN-TTS: Improving Zero-Shot Cross-Lingual Transfer for Multi-Lingual TTS with Style-Enhanced Normalization Diffusion Framework.....	4464
<i>Chung Tran, Chi Mai Luong, Sakriani Sakti</i>	

SEARCH METHODS AND DECODING ALGORITHMS FOR ASR

Average Token Delay: A Latency Metric for Simultaneous Translation	4469
<i>Yasumasa Kano, Katsuhito Sudoh, Satoshi Nakamura</i>	
Automatic Speech Recognition Transformer with Global Contextual Information Decoder.....	4474
<i>Yukun Qian, Xuyi Zhuang, Mingjiang Wang</i>	
Time-Synchronous One-Pass Beam Search for Parallel Online and Offline Transducers with Dynamic Block Training	4479
<i>Yui Sudo, Shakeel Muhammad, Yifan Peng, Shinji Watanabe</i>	
Prefix Search Decoding for RNN Transducers.....	4484
<i>Kiran Praveen, Advait Vinay Dhopeswarkar, Abhishek Pandey, Balaji Radhakrishnan</i>	
WhisperX: Time-Accurate Speech Transcription of Long-Form Audio.....	4489
<i>Max Bain, Jaesung Huh, Tengda Han, Andrew Zisserman</i>	
Implementing Contextual Biasing in GPU Decoder for Online ASR.....	4494
<i>Iuliia Nigmatulina, Srikanth Madikeri, Esau Villatoro-Tello, Petr Motlicek, Juan Zuluaga-Gomez, Karthik Pandia, Aravind Ganapathiraju</i>	

SPEECH SIGNAL ANALYSIS

MF-PAM: Accurate Pitch Estimation Through Periodicity Analysis and Multi-Level Feature Fusion.....	4499
<i>Woo-Jin Chung, Doyeon Kim, Soo-Whan Chung, Hong-Goo Kang</i>	
Enhancing Speech Articulation Analysis Using a Geometric Transformation of the X-Ray Microbeam Dataset.....	4504
<i>Ahmed Adel Attia, Mark Tiede, Carol Espy-Wilson</i>	
Matching Acoustic and Perceptual Measures of Phonation Assessment in Disordered Speech - A Case Study.....	4508
<i>Melanie Jouaiti, Pippa Kirby, Ravi Vaidyanathan</i>	
Improved Contextualized Speech Representations for Tonal Analysis	4513
<i>Jiahong Yuan, Xingyu Cai, Kenneth Church</i>	
A Study on the Importance of Formant Transitions for Stop-Consonant Classification in VCV Sequence.....	4518
<i>Siddarth Chandrasekar, Arvind Ramesh, Tilak Purohit, Prasanta Kumar Ghosh</i>	
FusedF0: Improving DNN-Based F0 Estimation by Fusion of Summary-Correlograms and Raw Waveform Representations of Speech Signals	4523
<i>Eray Eren, Lee Ngee Tan, Abeer Alwan</i>	

SPEECH EMOTION RECOGNITION 3

Improving Joint Speech and Emotion Recognition Using Global Style Tokens	4528
<i>Jehyun Kyung, Ju-Seok Seong, Jeong-Hwan Choi, Ye-Rin Jeoung, Joon-Hyuk Chang</i>	
Speech Emotion Recognition by Estimating Emotional Label Sequences with Phoneme Class Attribute	4533
<i>Ryotaro Nagase, Takahiro Fukumori, Yoichi Yamashita</i>	
Unsupervised Transfer Components Learning for Cross-Domain Speech Emotion Recognition	4538
<i>Shenjie Jiang, Peng Song, Shaokai Li, Keke Zhao, Wenming Zheng</i>	
Dual Memory Fusion for Multimodal Speech Emotion Recognition	4543
<i>Darshana Prisyad, Tharindu Fernando, Sridha Sridharan, Simon Denman, Clinton Fookes</i>	
Hybrid Dataset for Speech Emotion Recognition in Russian Language	4548
<i>Vladimir Kondratenko, Nikolay Karpov, Artem Sokolov, Nikita Savushkin, Oleg Kutuzov, Fyodor Minkin</i>	
Speech Emotion Recognition Using Decomposed Speech Via Multi-Task Learning	4553
<i>Jia-Hao Hsu, Chung-Hsien Wu, Yu-Hung Wei</i>	

CONNECTING SPEECH-SCIENCE AND SPEECH-TECHNOLOGY FOR CHILDREN'S SPEECH

Prospective Validation of Motor-Based Intervention with Automated Mispronunciation Detection of Rhotics in Residual Speech Sound Disorders	4558
<i>Nina R Benway, Jonathan L Preston</i>	
Classifying Rhoticity of /ɹ/ in Speech Sound Disorder Using Age-And-Sex Normalized Formants.....	4563
<i>Nina R Benway, Jonathan L Preston, Asif Salekin, Yi Xiao, Harshit Sharma, Tara McAllister</i>	
Acoustic-To-Articulatory Speech Inversion Features for Mispronunciation Detection of /ɹ/ in Child Speech Sound Disorders.....	4568
<i>Nina R Benway, Yashish M Siriwardena, Jonathan L Preston, Elaine Hitchcock, Tara McAllister, Carol Espy-Wilson</i>	
Using Commercial ASR Solutions to Assess Reading Skills in Children: A Case Report	4573
<i>Timothy Piton, Enno Hermann, Angela Pasqualotto, Marjolaine Cohen, Mathew Magimai.-Doss, Daphné Bavelier</i>	
Exploiting Diversity of Automatic Transcripts from Distinct Speech Recognition Techniques for Children's Speech.....	4578
<i>Christopher Gebauer, Lars Rumberg, Hanna Ehlert, Ulrike Lüdtkke, Joern Ostermann</i>	
Uncertainty Estimation for Connectionist Temporal Classification Based Automatic Speech Recognition	4583
<i>Lars Rumberg, Christopher Gebauer, Hanna Ehlert, Maren Wallbaum, Ulrike Lüdtkke, Joern Ostermann</i>	
BabySLM: Language-Acquisition-Friendly Benchmark of Self-Supervised Spoken Language Models.....	4588
<i>Marvin Lavechin, Yaya Sy, Hadrien Titeux, María Andrea Cruz Blandón, Okko Räsänen, Hervé Bredin, Emmanuel Dupoux, Alejandrina Cristia</i>	

Data Augmentation for Children ASR and Child-Adult Speaker Classification Using Voice Conversion Methods.....	4593
<i>Shuyang Zhao, Mittul Singh, Abraham Woubie, Reima Karhila</i>	
Developmental Articulatory and Acoustic Features for Six to Ten Year Old Children.....	4598
<i>Vishwas M. Shetty, Steven M. Lulich, Abeer Alwan</i>	
Automatically Predicting Perceived Conversation Quality in a Pediatric Sample Enriched for Autism.....	4603
<i>Yahan Yang, Sunghye Cho, Maxine Covello, Azia Knox, Osbert Bastani, James Weimer, Edgar Dobriban, Robert Schultz, Insup Lee, Julia Parish-Morris</i>	
An Equitable Framework for Automatically Assessing Children's Oral Narrative Language Abilities.....	4608
<i>Alexander Johnson, Hariram Veeramani, Natarajan Balaji Shankar, Abeer Alwan</i>	
An Analysis of Goodness of Pronunciation for Child Speech.....	4613
<i>Xinwei Cao, Zijian Fan, Torbjørn Svendsen, Giampiero Salvi</i>	
Measuring Language Development from Child-Centered Recordings.....	4618
<i>Yaya Sy, William N. Havard, Marvin Lavechin, Emmanuel Dupoux, Alejandrina Cristia</i>	
Speaking Clearly, Understanding Better: Predicting the L2 Narrative Comprehension of Chinese Bilingual Kindergarten Children Based on Speech Intelligibility Using a Machine Learning Approach.....	4623
<i>Hiuching Hung, Paula A. Pérez-Toro, Tomás Arias-Vergara, Andreas Maier, Elmar Nöth</i>	
Speech Breathing Behavior During Pauses in Children.....	4628
<i>Delphine Charuau, Béatrice Vaxelaire, Rudolph Sock</i>	
Understanding Spoken Language Development of Children with ASD Using Pre-Trained Speech Embeddings.....	4633
<i>Anfeng Xu, Rajat Hebbar, Rimita Lahiri, Tiantian Feng, Lindsay Butler, Lue Shen, Helen Tager-Flusberg, Shrikanth Narayanan</i>	
Measuring Phonological Precision in Children with Cleft Lip and Palate.....	4638
<i>Tomás Arias-Vergara, Elizabeth Londoño-Mora, Paula A. Pérez-Toro, Maria Schuster, Elmar Nöth, Juan Rafael Orozco-Arroyave, Andreas Maier</i>	
A Study on Using Duration and Formant Features in Automatic Detection of Speech Sound Disorder in Children.....	4643
<i>Si-Ioi Ng, Cymie Wing-Yee Ng, Tan Lee</i>	
Influence of Utterance and Speaker Characteristics on the Classification of Children with Cleft Lip and Palate.....	4648
<i>Ilja Baumann, Dominik Wagner, Franziska Braun, Sebastian P. Bayerl, Elmar Nöth, Korbinian Riedhammer, Tobias Bocklet</i>	

DIALOG MANAGEMENT

Parameter-Efficient Low-Resource Dialogue State Tracking by Prompt Tuning.....	4653
<i>Mingyu Derek Ma, Jiun-Yu Kao, Shuyang Gao, Arpit Gupta, Di Jin, Tagyoung Chung, Nanyun Peng</i>	
An Autoregressive Conversational Dynamics Model for Dialogue Systems.....	4658
<i>Matthew McNeill, Rivka Levitan</i>	

Style-Transfer Based Speech and Audio-Visual Scene Understanding for Robot Action Sequence Acquisition from Videos.....	4663
<i>Chiori Hori, Puyuan Peng, David Harwath, Xinyu Liu, Kei Ota, Siddarth Jain, Radu Corcodel, Devesh Jha, Diego Romeres, Jonathan Le Roux</i>	
Speech Aware Dialog System Technology Challenge (DSTC11).....	4668
<i>Hagen Soltau, Izhak Shafran, Mingqiu Wang, Abhinav Rastogi, Jeffrey Zhao, Ye Jia, Wei Han, Yuan Cao, Aramys Miranda</i>	
Knowledge-Retrieval Task-Oriented Dialog Systems with Semi-Supervision.....	4673
<i>Yucheng Cai, Hong Liu, Zhijian Ou, Yi Huang, Junlan Feng</i>	
Tracking Must Go on : Dialogue State Tracking with Verified Self-Training	4678
<i>Jihyun Lee, Chaebin Lee, Yunsu Kim, Gary Geunbae Lee</i>	

SPEAKER RECOGNITION 2

Ordered and Binary Speaker Embedding	4683
<i>Jiaying Wang, Xianglong Wang, Namin Wang, Lantian Li, Dong Wang</i>	
Self-FiLM: Conditioning GANs with Self-Supervised Representations for Bandwidth Extension Based Speaker Recognition	4688
<i>Saurabh Kataria, Jesús Villalba, Laureano Moro-Velazquez, Thomas Thebaud, Najim Dehak</i>	
Curriculum Learning for Self-Supervised Speaker Verification.....	4693
<i>Hee-Soo Heo, Jee-Weon Jung, Jingu Kang, Young-Ki Kwon, Bong-Jin Lee, You Jin Kim, Joon Son Chung</i>	
Introducing Self-Supervised Phonetic Information for Text-Independent Speaker Verification	4698
<i>Ziyang Zhang, Wu Guo, Bin Gu</i>	
A Teacher-Student Approach for Extracting Informative Speaker Embeddings from Speech Mixtures	4703
<i>Tobias Cord-Landwehr, Christoph Boeddeker, Catalin Zorila, Rama Doddipatla, Reinhold Haeb-Umbach</i>	
Experimenting with Additive Margins for Contrastive Self-Supervised Speaker Verification.....	4708
<i>Theo Lepage, Reda Dehak</i>	

PHONETICS, PHONOLOGY, AND PROSODY 2

Nonbinary American English Speakers Encode Gender in Vowel Acoustics.....	4713
<i>Maxwell Hope, Charlotte Ward, Jason Lilley</i>	
Coarticulation of Sibe Vowels and Dorsal Fricatives in Spontaneous Speech: An Acoustic Study.....	4718
<i>Jared Sharp, Matthew Faytak, Hasutai Fei Xiong Liu</i>	
Using Speech Synthesis to Explain Automatic Speaker Recognition: A New Application of Synthetic Speech	4723
<i>Georgina Brown, Christin Kirchhübel, Ramiz Cuthbert</i>	
Same F0, Different Tones: A Multidimensional Investigation of Zhangzhou Tones	4728
<i>Yishan Huang</i>	

Discovering Phonetic Feature Event Patterns in Transformer Embeddings	4733
<i>Patrick Cormac English, John D. Kelleher, Julie Carson-Berndsen</i>	
A System for Generating Voice Source Signals that Implements the Transformed LF-Model Parameter Control.....	4738
<i>Zihan Wang, Christer Gobl</i>	
Speaker-Independent Speech Inversion for Estimation of Nasalance	4743
<i>Yashish M Siriwardena, Carol Espy-Wilson, Suzanne Boyce, Mark Tiede, Liran Oren</i>	
Effects of Tonal Coarticulation and Prosodic Positions on Tonal Contours of Low Rising Tones: In the Case of Xiamen Dialect.....	4748
<i>Yiying Hu, Hui Feng, Qinghua Zhao, Aijun Li</i>	
Durational and Non-Durational Correlates of Lexical and Derived Geminates in Arabic	4753
<i>Amel Issa</i>	
Mapping Phonemes to Acoustic Symbols and Codes Using Synchrony in Speech Modulation Vectors Estimated by the Travellingwave Filter Bank.....	4758
<i>Ashwin Rao</i>	
Rhythmic Characteristics of L2 German Speech by Advanced Chinese Learners	4763
<i>Lindun Ge, Min Xu, Hongwei Ding</i>	
(Dis)agreement and Preference Structure Are Reflected in Matching Along Distinct Acoustic- Prosodic Features	4768
<i>Anneliese Kelterer, Margaret Zellers, Barbara Schuppler</i>	
Vowel Reduction by Greek-Speaking Children: The Effect of Stress and Word Length.....	4773
<i>Polychronia Christodoulidou, Katerina Nicolaidis, Dimitrios Stamovlasis</i>	
Pitch Distributions in a Very Large Corpus of Spontaneous Finnish Speech	4778
<i>Mietta Lennes, Minnaleena Toivola</i>	
Speech Enhancement Patterns in Human-Robot Interaction: A Cross-Linguistic Perspective.....	4783
<i>Jacek Kudera, Katharina Zahner-Ritter, Jakob Engel, Nathalie Elsässer, Philipp Hutmacher, Carolin Worstbrock</i>	

SPEECH SYNTHESIS: EXPRESSIVITY

Controllable Generation of Artificial Speaker Embeddings Through Discovery of Principal Directions	4788
<i>Florian Lux, Pascal Tilli, Sarina Meyer, Ngoc Thang Vu</i>	
Dual Audio Encoders Based Mandarin Prosodic Boundary Prediction by Using Multi-Granularity Prosodic Representations.....	4793
<i>Ruishan Li, Yingming Gao, Yanlu Xie, Dengfeng Ke, Jinsong Zhang</i>	
NoreSpeech: Knowledge Distillation Based Conditional Diffusion Model for Noise-Robust Expressive TTS	4798
<i>Dongchao Yang, Songxiang Liu, Helin Wang, Jianwei Yu, Chao Weng, Yuexian Zou</i>	
MaskedSpeech: Context-Aware Speech Synthesis with Masking Strategy.....	4803
<i>Ya-Jie Zhang, Wei Song, Yanghao Yue, Zhengchen Zhang, Youzheng Wu, Xiaodong He</i>	

Narrator Or Character: Voice Modulation in an Expressive Multi-Speaker TTS	4808
<i>Tankala Pavan Kalyan, Preeti Rao, Preethi Jyothi, Pushpak Bhattacharyya</i>	
CASEIN: Cascading Explicit and Implicit Control for Fine-Grained Emotion Intensity Regulation	4813
<i>Yuhao Cui, Xiongwei Wang, Zhongzhou Zhao, Wei Zhou, Haiqing Chen</i>	
Semi-Supervised Learning for Continuous Emotional Intensity Controllable Speech Synthesis with Disentangled Representations.....	4818
<i>Yoori Oh, Juheon Lee, Yoseob Han, Kyogu Lee</i>	
Expresso: A Benchmark and Analysis of Discrete Expressive Speech Resynthesis.....	4823
<i>Tu Anh Nguyen, Wei-Ning Hsu, Antony D'Avirro, Bowen Shi, Itai Gat, Maryam Fazel-Zarani, Tal Remez, Jade Copet, Gabriel Synnaeve, Michael Hassid, Felix Kreuk, Yossi Adi, Emmanuel Dupoux</i>	
ComedicSpeech: Text to Speech for Stand-Up Comedies in Low-Resource Scenarios.....	4828
<i>Yuyue Wang, Huan Xiao, Yihan Wu, Ruihua Song</i>	
Neural Speech Synthesis with Enriched Phrase Boundaries	4833
<i>Marie Kunešová, Jindrich Matoušek</i>	
Cross-Lingual Prosody Transfer for Expressive Machine Dubbing	4838
<i>Jakub Swiatkowski, Duo Wang, Mikolaj Babianski, Patrick Lumban Tobing, Ravichander Vipperla, Vincent Pollet</i>	
Synthesis After a Couple PINTs: Investigating the Role of Pause-Internal Phonetic Particles in Speech Synthesis and Perception	4843
<i>Mikey Elmers, Johannah O'Mahony, Éva Székely</i>	
Accentor: An Explicit Lexical Stress Model for TTS Systems	4848
<i>Diana Geneva, Georgi Shopov, Kostadin Garov, Maria Todorova, Stefan Gerdjikov, Stoyan Mihov</i>	
A Neural TTS System with Parallel Prosody Transfer from Unseen Speakers	4853
<i>Slava Shechtman, Raul Fernandez</i>	
Diverse and Expressive Speech Prosody Prediction with Denoising Diffusion Probabilistic Model.....	4858
<i>Xiang Li, Songxiang Liu, Max W. Y. Lam, Zhiyong Wu, Chao Weng, Helen Meng</i>	
Prosody Modeling with 3D Visual Information for Expressive Video Dubbing	4863
<i>Zhihan Yang, Shansong Liu, Xu Li, Haozhe Wu, Zhiyong Wu, Ying Shan, Jia Jia</i>	
LightClone: Speaker-Guided Parallel Subnet Selection for Few-Shot Voice Cloning	4868
<i>Jie Wu, Jian Luan, Yujun Wang</i>	

VOLUME 8

EE-TTS: Emphatic Expressive TTS with Linguistic Information.....	4873
<i>Yi Zhong, Chen Zhang, Xule Liu, Chenxi Sun, Weishan Deng, Haifeng Hu, Zhongqian Sun</i>	
Stochastic Pitch Prediction Improves the Diversity and Naturalness of Speech in Glow-TTS	4878
<i>Sewade Ogun, Vincent Colotte, Emmanuel Vincent</i>	
ContextSpeech: Expressive and Efficient Text-To-Speech for Paragraph Reading.....	4883
<i>Yujia Xiao, Shaofei Zhang, Xi Wang, Xu Tan, Lei He, Sheng Zhao, Frank K. Soong, Tan Lee</i>	

PromptStyle: Controllable Style Transfer for Text-To-Speech with Natural Language Descriptions	4888
<i>Guanghou Liu, Yongmao Zhang, Yi Lei, Yunlin Chen, Rui Wang, Lei Xie, Zhifei Li</i>	
Creating Personalized Synthetic Voices from Post-Glossectomy Speech with Guided Diffusion Models.....	4893
<i>Yusheng Tian, Guangyan Zhang, Tan Lee</i>	

SPEECH RECOGNITION: SIGNAL PROCESSING, ACOUSTIC MODELING, ROBUSTNESS, ADAPTATION 5

Towards Multi-Task Learning of Speech and Speaker Recognition.....	4898
<i>Nik Vaessen, David A. Van Leeuwen</i>	
Regarding Topology and Variant Frame Rates for Differentiable WFST-Based End-To-End ASR	4903
<i>Zeyu Zhao, Peter Bell</i>	
2-Bit Conformer Quantization for Automatic Speech Recognition.....	4908
<i>Oleg Rybakov, Phoenix Meadowlark, Shaojin Ding, David Qiu, Jian Li, David Rim, Yanzhang He</i>	
Time-Domain Speech Enhancement for Robust Automatic Speech Recognition	4913
<i>Yufeng Yang, Ashutosh Pandey, Deliang Wang</i>	
Multi-Channel Multi-Speaker Transformer for Speech Recognition	4918
<i>Guo Yifan, Tian Yao, Suo Hongbin, Wan Yulong</i>	
Fake the Real: Backdoor Attack on Deep Speech Classification Via Voice Conversion.....	4923
<i>Zhe Ye, Terui Mao, Li Dong, Diqun Yan</i>	
Dialect Speech Recognition Modeling Using Corpus of Japanese Dialects and Self-Supervised Learning-Based Model XLSR.....	4928
<i>Shogo Miwa, Atsuhiko Kai</i>	
Contextualized End-To-End Speech Recognition with Contextual Phrase Prediction Network	4933
<i>Kaixun Huang, Ao Zhang, Zhanheng Yang, Pengcheng Guo, Bingshen Mu, Tianyi Xu, Lei Xie</i>	
Competitive and Resource Efficient Factored Hybrid HMM Systems Are Simpler than You Think.....	4938
<i>Tina Raissi, Christoph Lüscher, Moritz Gunz, Ralf Schlüter, Hermann Ney</i>	
MMSpeech: Multi-Modal Multi-Task Encoder-Decoder Pre-Training for Speech Recognition.....	4943
<i>Xiaohuan Zhou, Jiaming Wang, Zeyu Cui, Shiliang Zhang, Zhijie Yan, Jingren Zhou, Chang Zhou</i>	
Biased Self-Supervised Learning for ASR	4948
<i>Florian L. Kreyssig, Yangyang Shi, Jinxi Guo, Leda Sari, Abdel-Rahman Mohamed, Philip C. Woodland</i>	
A Unified Recognition and Correction Model Under Noisy and Accent Speech Conditions.....	4953
<i>Zhao Yang, Dianwen Ng, Chong Zhang, Rui Jiang, Wei Xi, Yukun Ma, Chongjia Ni, Jizhong Zhao, Bin Ma, Eng Siong Chng</i>	
Wav2vec 2.0 ASR for Cantonese-Speaking Older Adults in a Clinical Setting	4958
<i>Ranzo Huang, Brian Mak</i>	
BAT: Boundary Aware Transducer for Memory-Efficient and Low-Latency ASR	4963
<i>Keyu An, Xian Shi, Shiliang Zhang</i>	

Bayes Risk Transducer: Transducer with Controllable Alignment Prediction.....	4968
<i>Jinchuan Tian, Jianwei Yu, Hangting Chen, Brian Yan, Chao Weng, Dong Yu, Shinji Watanabe</i>	
Multi-View Frequency-Attention Alternative to CNN Frontends for Automatic Speech Recognition	4973
<i>Belen Alastruey, Lukas Drude, Jahn Heymann, Simon Wiesler</i>	

SPEECH, VOICE, AND HEARING DISORDERS 2

Investigating the Dynamics of Hand and Lips in French Cued Speech Using Attention Mechanisms and CTC-Based Decoding	4978
<i>Sanjana Sankar, Denis Beautemps, Frédéric Elisei, Olivier Perrotin, Thomas Hueber</i>	
Hearing Loss Affects Emotion Perception in Older Adults: Evidence from a Prosody-Semantics Stroop Task	4983
<i>Yingyang Wang, Min Xu, Jing Shao, Lan Wang, Nan Yan</i>	
Cochlear-Implant Listeners Listening to Cochlear-Implant Simulated Speech.....	4988
<i>Fanhui Kong, Nengheng Zheng, Xianren Wang, Hao He, Jan W. H. Schnupp, Qinglin Meng</i>	
Validation of a Task-Independent Cepstral Peak Prominence Measure with Voice Activity Detection	4993
<i>Olivia M. Murton, Abigail E. Haenssler, Marc F. Maffei, Kathryn P. Connaghan, Jordan Green</i>	
Score-Balanced Loss for Multi-Aspect Pronunciation Assessment.....	4998
<i>Heejin Do, Yunsu Kim, Gary Geunbae Lee</i>	
Federated Learning for Secure Development of AI Models for Parkinson’s Disease Detection Using Speech from Different Languages	5003
<i>Soroosh Tayebi Arasteh, Cristian David Rios-Urrego, Elmar Nöth, Andreas Maier, Seung Hee Yang, Jan Rusz, Juan Rafael Orozco-Arroyave</i>	
F0inTFS: A Lightweight Periodicity Enhancement Strategy for Cochlear Implants.....	5008
<i>Huali Zhou, Fanhui Kong, Nengheng Zheng, Qinglin Meng</i>	
Differentiating Acoustic and Physiological Features in Speech for Hypoxia Detection	5013
<i>Benjamin O'Brien, Adrien Gresse, Jean-Baptiste Billaud, Guilhem Belda, Jean-François Bonastre</i>	
Mandarin Electrolaryngeal Speech Voice Conversion Using Cross-Domain Features	5018
<i>Hsin-Hao Chen, Yung-Lun Chien, Ming-Chi Yen, Shu-Wei Tsai, Tai-Shih Chi, Hsin-Min Wang, Yu Tsao</i>	
Audio-Visual Mandarin Electrolaryngeal Speech Voice Conversion	5023
<i>Yung-Lun Chien, Hsin-Hao Chen, Ming-Chi Yen, Shu-Wei Tsai, Hsin-Min Wang, Yu Tsao, Tai-Shih Chi</i>	
Which Aspects of Motor Speech Disorder Are Captured by Mel Frequency Cepstral Coefficients? Evidence from the Change in STN-DBS Conditions in Parkinson’s Disease	5027
<i>Vojtech Illner, Petr Krýže, Jan Švihlík, Mário Sousa, Paul Krack, Elina Tripoliti, Robert Jech, Jan Rusz</i>	
Detecting Manifest Huntington's Disease Using Vocal Data.....	5032
<i>Vinod Subramanian, Namhee Kwon, Raymond Brueckner, Nate Blaylock, Henry O'Connell, Luis Sierra, Clementina Ullman, Karen Hildebrand, Simon Laganieri</i>	

Exploring Multi-Task Learning and Data Augmentation in Dementia Detection with Self-Supervised Pretrained Models	5037
<i>Minchuan Chen, Chenfeng Miao, Jun Ma, Shaojun Wang, Jing Xiao</i>	

SPEECH ACTIVITY DETECTION AND MODELING

GL-SSD: Global and Local Speech Style Disentanglement by Vector Quantization for Robust Sentence Boundary Detection in Speech Stream.....	5042
<i>Kuncai Zhang, Wei Zhou, Pengcheng Zhu, Haiqing Chen</i>	
Semantic VAD: Low-Latency Voice Activity Detection for Speech Interaction	5047
<i>Mohan Shi, Yuchun Shu, Lingyun Zuo, Qian Chen, Shiliang Zhang, Jie Zhang, Li-Rong Dai</i>	
Dynamic Encoder RNN for Online Voice Activity Detection in Adverse Noise Conditions	5052
<i>Prithvi R. R. Gudepu, Jayesh M. Koroth, Kamini Sabu, Mahaboob Ali Basha Shaik</i>	
Point to the Hidden: Exposing Speech Audio Splicing Via Signal Pointer Nets	5057
<i>Denise Moussa, Germans Hirsch, Sebastian Wankerl, Christian Riess</i>	
Real-Time Causal Spectro-Temporal Voice Activity Detection Based on Convolutional Encoding and Residual Decoding.....	5062
<i>Jingyuan Wang, Jie Zhang, Li-Rong Dai</i>	
SVVAD: Personal Voice Activity Detection for Speaker Verification.....	5067
<i>Zuheng Kang, Jianzong Wang, Junqing Peng, Jing Xiao</i>	

MULTILINGUAL MODELS FOR ASR

Learning Cross-Lingual Mappings for Data Augmentation to Improve Low-Resource Speech Recognition	5072
<i>Muhammad Umar Farooq, Thomas Hain</i>	
AfriNames: Most ASR Models "Butcher" African Names.....	5077
<i>Tobi Olatunji, Tejumade Afonja, Bonaventure F. P. Dossou, Atnafu Lambebo Tonja, Chris Chinenye Emezue, Amina Mardiyah Rufai, Sahib Singh</i>	
Towards Dialect-Inclusive Recognition in a Low-Resource Language: Are Balanced Corpora the Answer?.....	5082
<i>Liam Lonergan, Mengjie Qian, Neasa Ní Chiaráin, Christer Gobl, Ailbhe Ní Chasaide</i>	
Svarah: Evaluating English ASR Systems on Indian Accents	5087
<i>Tahir Javed, Sakshi Joshi, Vignesh Nagarajan, Sai Sundaresan, Janki Nawale, Abhigyan Raman, Kaushal Bhogale, Pratyush Kumar, Mitesh M. Khapra</i>	
N-Shot Benchmarking of Whisper on Diverse Arabic Speech Recognition.....	5092
<i>Bashar Talafha, Abdul Waheed, Muhammad Abdul-Mageed</i>	
The MALACH Corpus: Results with End-To-End Architectures and Pretraining	5097
<i>Michael Picheny, Qin Yang, Daiheng Zhang, Lining Zhang</i>	

SPEECH ENHANCEMENT AND BANDWIDTH EXPANSION

Unsupervised Speech Enhancement with Deep Dynamical Generative Speech and Noise Models	5102
<i>Xiaoyu Lin, Simon Leglaive, Laurent Girin, Xavier Alameda-Pineda</i>	

Noise-Robust Bandwidth Expansion for 8K Speech Recordings.....	5107
<i>Yin-Tse Lin, Bo-Hao Su, Chi-Han Lin, Shih-Chan Kuo, Jyh-Shing Roger Jang, Chi-Chun Lee</i>	
MdctGAN: Taming Transformer-Based GAN for Speech Super-Resolution with Modified DCT Spectra.....	5112
<i>Chenhao Shuai, Chaohua Shi, Lu Gan, Hongqing Liu</i>	
Zoneformer: On-Device Neural Beamformer for In-Car Multi-Zone Speech Separation, Enhancement and Echo Cancellation	5117
<i>Yong Xu, Vinay Kothapally, Meng Yu, Shixiong Zhang, Dong Yu</i>	
Low-Complexity Broadband Beampattern Synthesis Using Array Response Control.....	5122
<i>Jiayi Xu, Jian Li, Weixin Meng, Xiaodong Li, Chengshi Zheng</i>	
A GAN Speech Inpainting Model for Audio Editing Software	5127
<i>Haixin Zhao</i>	

ARTICULATION

Deep Speech Synthesis from MRI-Based Articulatory Representations.....	5132
<i>Peter Wu, Tingle Li, Yijing Lu, Yubin Zhang, Jiachen Lian, Alan W Black, Louis Goldstein, Shinji Watanabe, Gopala K. Anumanchipalli</i>	
Learning to Compute the Articulatory Representations of Speech with the MIRRORNET.....	5137
<i>Yashish M Siriwardena, Carol Espy-Wilson, Shihab Shamma</i>	
Generating High-Resolution 3D Real-Time MRI of the Vocal Tract	5142
<i>Martin Strauch, Antoine Serrurier</i>	
Exploring a Classification Approach Using Quantised Articulatory Movements for Acoustic to Articulatory Inversion	5147
<i>Jesuraj Bandekar, Sathvik Udupa, Prasanta Kumar Ghosh</i>	

NEURAL PROCESSING OF SPEECH AND LANGUAGE: ENCODING AND DECODING THE DIVERSE AUDITORY BRAIN

MEG Encoding Using Word Context Semantics in Listening Stories.....	5152
<i>Subba Reddy Oota, Nathan Trouvain, Frederic Alexandre, Xavier Hinaut</i>	
Investigating the Cortical Tracking of Speech and Music with Sung Speech	5157
<i>Giorgia Cantisani, Amirhossein Chalehchaleh, Giovanni Di Liberto, Shihab Shamma</i>	
Coherence Estimation Tracks Auditory Attention in Listeners with Hearing Impairment	5162
<i>Oskar Keding, Emina Alickovic, Martin A. Skoglund, Maria Sandsten</i>	
Speech Taskonomy: Which Speech Tasks Are the Most Predictive of fMRI Brain Activity?.....	5167
<i>Subba Reddy Oota, Veeral Agarwal, Mounika Marreddy, Manish Gupta, Raju Bapi</i>	
Exploring Auditory Attention Decoding Using Speaker Features	5172
<i>Zelin Qiu, Jianjun Gu, Dingding Yao, Junfeng Li</i>	
Enhancing the EEG Speech Match Mismatch Tasks with Word Boundaries	5177
<i>Akshara Soman, Vidhi Sinha, Sriram Ganapathy</i>	

Similar Hierarchical Representation of Speech and Other Complex Sounds in the Brain and Deep Residual Networks: An MEG Study.....	5182
<i>Tzu-Han Zoe Cheng, Kuan-Lin Chen, Juliane Schubert, Ya-Ping Chen, Tim Brown, John Iversen</i>	
Effects of Spectral Degradation on the Cortical Tracking of the Speech Envelope.....	5187
<i>Alexis Deighton Macintyre, Tobias Goehring</i>	
Effects of Spectral and Temporal Modulation Degradation on Intelligibility and Cortical Tracking of Speech Signals	5192
<i>Ignacio Calderon De Palma, Laura S. Lopez, Alejandro Lopez Valdes</i>	

PERCEPTION OF PARALINGUISTICS

Transfer Learning for Personality Perception Via Speech Emotion Recognition.....	5197
<i>Yuanchao Li, Peter Bell, Catherine Lai</i>	
A Stimulus-Organism-Response Model of Willingness to Buy from Advertising Speech Using Voice Quality	5202
<i>Mizuki Nagano, Yusuke Ijima, Sadao Hiroya</i>	
Voice Passing : A Non-Binary Voice Gender Prediction System for Evaluating Transgender Voice Transition.....	5207
<i>David Doukhan, Simon Devauchelle, Lucile Girard-Monneron, Mía Chávez Ruz, V. Chaddouk, Isabelle Wagner, Albert Rilliard</i>	
Influence of Personal Traits on Impressions of One's Own Voice	5212
<i>Hikaru Yanagida, Yusuke Ijima, Naohiro Tawara</i>	
Pardon My Disfluency: The Impact of Disfluency Effects on the Perception of Speaker Competence and Confidence	5217
<i>Ambika Kirkland, Joakim Gustafson, Éva Székely</i>	
Cross-Linguistic Emotion Perception in Human and TTS Voices.....	5222
<i>Iona Gessinger, Michelle Cohn, Benjamin R. Cowan, Georgia Zellou, Bernd Möbius</i>	

TECHNOLOGIES FOR CHILD SPEECH PROCESSING

Joint Learning Feature and Model Adaptation for Unsupervised Acoustic Modelling of Child Speech	5227
<i>Richeng Duan</i>	
Automatic Assessment of Oral Reading Accuracy for Reading Diagnostics	5232
<i>Bo Molenaar, Cristian Tejedor-Garcia, Catia Cucchiarini, Helmer Strik</i>	
An ASR-Enabled Reading Tutor: Investigating Feedback to Optimize Interaction for Learning to Read.....	5237
<i>Yu Bai, Ferdy Hubers, Catia Cucchiarini, Roeland Van Hout, Helmer Strik</i>	
Adaptation of Whisper Models to Child Speech Recognition	5242
<i>Rishabh Jain, Andrei Barcovschi, Mariam Yiwere, Peter Corcoran, Horia Cucu</i>	

SHOW AND TELL: MEDIA AND COMMERCIAL APPLICATIONS

Let's Give a Voice to Conversational Agents in Virtual Reality	5247
<i>Michele Yin, Gabriel Roccabruna, Abhinav Azad, Giuseppe Riccardi</i>	
FOOCTTS: Generating Arabic Speech with Acoustic Environment for Football Commentator	5249
<i>Massa Baali, Ahmed M. Ali</i>	
Video Summarization Leveraging Multimodal Information for Presentations.....	5251
<i>Hanchao Liu, Dapeng Chen, Rongjun Li, Wenyuan Xue, Wei Peng</i>	
What Questions Are My Customers Asking?: Towards Actionable Insights from Customer Questions in Contact Center Calls.....	5253
<i>Varun Nathan, Devashish Deshpande, Ayush Kumar, Cijo George, Jithendra Vepa</i>	
COnVoy: A Contact Center Operated Pipeline for Voice of Customer Discovery.....	5255
<i>Rishabh Tripathi, Digvijay Anil Ingle, Ayush Kumar, Cijo George, Jithendra Vepa</i>	
NeMo Forced Aligner and Its Application to Word Alignment for Subtitle Generation.....	5257
<i>Elena Rastorgueva, Vitaly Lavrukhin, Boris Ginsburg</i>	
CauSE: Causal Search Engine for Understanding Contact-Center Conversations	5259
<i>Anup Pattnaik, Tanay Narshana, Aashraya Sachdeva, Cijo George, Jithendra Vepa</i>	
Tailored Real-Time Call Summarization System for Contact Centers	5261
<i>Aashraya Sachdeva, Sai Nishanth Padala, Anup Pattnaik, Varun Nathan, Cijo George, Ayush Kumar, Jithendra Vepa</i>	
Federated Learning Toolkit with Voice-Based User Verification Demo.....	5263
<i>Prathamesh Mandke, Rachel Oberst, Matthias Reisser, Av?it Chakraborty, Christos Louizos, Joseph Soriaga, Daniel Madrigal, Andre Manoel, Nalin Singal, Jeff Omhover, Robert Sim</i>	
Learning When to Speak: Latency and Quality Trade-Offs for Simultaneous Speech-To-Speech Translation with Offline Models.....	5265
<i>Liam Dugan, Anshul Wadhawan, Kyle Spence, Chris Callison-Burch, Morgan McGuire, Victor Zordan</i>	
Fast Enrollable Streaming Keyword Spotting System: Training and Inference Using a Web Browser	5267
<i>Namhyun Cho, Sunmin Kim, Yoseb Kang, Heeman Kim</i>	
Cross-Lingual/Cross-Channel Intent Detection in Contact-Center Conversations.....	5269
<i>Suraj Agrawal, Aashraya Sachdeva, Soumya Jain, Cijo George, Jithendra Vepa</i>	

SPEAKER AND LANGUAGE IDENTIFICATION 3

One-Step Knowledge Distillation and Fine-Tuning in Using Large Pre-Trained Self-Supervised Learning Models for Speaker Verification.....	5271
<i>Jungwoo Heo, Chan-Yeong Lim, Ju-Ho Kim, Hyun-Seo Shin, Ha-Jin Yu</i>	
Defense Against Adversarial Attacks on Audio DeepFake Detection	5276
<i>Piotr Kawa, Marcin Plata, Piotr Syga</i>	
A Conformer-Based Classifier for Variable-Length Utterance Processing in Anti-Spoofing.....	5281
<i>Eros Rosello, Alejandro Gomez-Alanis, Angel M. Gomez, Antonio Peinado</i>	

Conformer-Based Language Embedding with Self-Knowledge Distillation for Spoken Language Identification	5286
<i>Feng Wang, Lingyan Huang, Tao Li, Qingyang Hong, Lin Li</i>	
CommonAccent: Exploring Large Acoustic Pretrained Models for Accent Classification Based on Common Voice	5291
<i>Juan Zuluaga-Gomez, Sara Ahmed, Danielius Visockas, Cem Subakan</i>	
From Adaptive Score Normalization to Adaptive Data Normalization for Speaker Verification Systems.....	5296
<i>Sandro Cumani, Salvatore Sarni</i>	
CAM++: A Fast and Efficient Network for Speaker Verification Using Context-Aware Masking	5301
<i>Hui Wang, Siqi Zheng, Yafeng Chen, Luyao Cheng, Qian Chen</i>	
North Sámi Dialect Identification with Self-Supervised Speech Models.....	5306
<i>Sofoklis Kakouros, Katri Hiovain-Asikainen</i>	
Encoder-Decoder Multimodal Speaker Change Detection.....	5311
<i>Jee-Weon Jung, Soonshin Seo, Hee-Soo Heo, Geonmin Kim, You Jin Kim, Young-Ki Kwon, Minjae Lee, Bong-Jin Lee</i>	
Disentangled Representation Learning for Multilingual Speaker Recognition	5316
<i>Kihyun Nam, Youkyum Kim, Jaesung Huh, Hee-Soo Heo, Jee-Weon Jung, Joon Son Chung</i>	
A Compact End-To-End Model with Local and Global Context for Spoken Language Identification	5321
<i>Fei Jia, Nithin Rao Koluguri, Jagadeesh Balam, Boris Ginsburg</i>	
On the Robustness of Arabic Speech Dialect Identification.....	5326
<i>Peter Sullivan, Abdelrahim Elmadany, Muhammad Abdul-Mageed</i>	
Adaptive Neural Network Quantization for Lightweight Speaker Verification.....	5331
<i>Haoyu Wang, Bei Liu, Yifei Wu, Yanmin Qian</i>	
Adversarial Diffusion Probability Model for Cross-Domain Speaker Verification Integrating Contrastive Loss	5336
<i>Xinmei Su, Xiang Xie, Fengrun Zhang, Chenguang Hu</i>	
Spoofing Attacker Also Benefits from Self-Supervised Pretrained Model.....	5346
<i>Aoi Ito, Shota Horiguchi</i>	
Label Aware Speech Representation Learning for Language Identification.....	5351
<i>Shikhar Vashishth, Shikhar Bharadwaj, Sriram Ganapathy, Ankur Bapna, Min Ma, Wei Han, Vera Axelrod, Partha Talukdar</i>	
Exploring the Impact of Back-End Network on Wav2vec 2.0 for Dialect Identification	5356
<i>Qibao Luo, Ruohua Zhou</i>	
Improving Speaker Verification with Self-Pretrained Transformer Models	5361
<i>Junyi Peng, Oldrich Plchot, Themis Stafylakis, Ladislav Mosner, Lukáš Burget, Jan "Honza" Cernocký</i>	
Handling the Alignment for Wake Word Detection: A Comparison Between Alignment-Based, Alignment-Free and Hybrid Approaches.....	5366
<i>Vinicius Ribeiro, Yiteng Huang, Yuan Shangguan, Zhaojun Yang, Li Wan, Ming Sun</i>	

ANALYSIS OF SPEECH AND AUDIO SIGNALS 4

What Do Self-Supervised Speech Representations Encode? An Analysis of Languages, Varieties, Speaking Styles and Speakers	5371
<i>Julian Linke, Mate Kadar, Gergely Dosinszky, Peter Mihajlik, Gernot Kubin, Barbara Schuppler</i>	
A Compressed Synthetic Speech Detection Method with Compression Feature Embedding	5376
<i>Jinghong Zhang, Xiaowei Yi, Xianfeng Zhao</i>	
Outlier-Aware Inlier Modeling and Multi-Scale Scoring for Anomalous Sound Detection Via Multitask Learning	5381
<i>Yucong Zhang, Suo Hongbin, Yulong Wan, Ming Li</i>	
MOSLight: A Lightweight Data-Efficient System for Non-Intrusive Speech Quality Assessment.....	5386
<i>Zitong Li, Wei Li</i>	
A Multi-Scale Attentive Transformer for Multi-Instrument Symbolic Music Generation.....	5391
<i>Xipin Wei, Junhui Chen, Zirui Zheng, Li Guo, Lantian Li, Dong Wang</i>	
MTANet: Multi-Band Time-Frequency Attention Network for Singing Melody Extraction from Polyphonic Music	5396
<i>Yuan Gao, Ying Hu, Liusong Wang, Hao Huang, Liang He</i>	
Xiaoicesing 2: A High-Fidelity Singing Voice Synthesizer Based on Generative Adversarial Network.....	5401
<i>Wang Chunhui, Chang Zeng, Xing He</i>	
Do Vocal Breath Sounds Encode Gender Cues for Automatic Gender Classification?	5406
<i>Mohammad Shaique Solanki, Ashutosh Bharadwaj, Jeevan Kylash, Prasanta Kumar Ghosh</i>	
Automatic Exploration of Optimal Data Processing Operations for Sound Data Augmentation Using Improved Differentiable Automatic Data Augmentation	5411
<i>Toki Sugiura, Hiromitsu Nishizaki</i>	
A Snoring Sound Dataset for Body Position Recognition: Collection, Annotation, and Analysis	5416
<i>Li Xiao, Xiuping Yang, Xinhong Li, Weiping Tu, Xiong Chen, Weiyan Yi, Jie Lin, Yuhong Yang, Yanzhen Ren</i>	
RMVPE: A Robust Model for Vocal Pitch Estimation in Polyphonic Music	5421
<i>Haojie Wei, Xueke Cao, Tangpeng Dan, Yueguo Chen</i>	
Spatialization Quality Metric for Binaural Speech.....	5426
<i>Pranay Manocha, Israel Dejene Gebru, Anurag Kumar, Dejan Markovic, Alexander Richard</i>	
AsthmaSCeLNet: A Lightweight Supervised Contrastive Embedding Learning Framework for Asthma Classification Using Lung Sounds	5431
<i>Arka Roy, Udit Satija</i>	
Patch-Mix Contrastive Learning with Audio Spectrogram Transformer on Respiratory Sound Classification	5436
<i>Sangmin Bae, June-Woo Kim, Won-Yang Cho, Hyerim Baek, Soyoun Son, Byungjo Lee, Changwan Ha, Kyongpil Tae, Sungnyun Kim, Se-Young Yun</i>	

Remote Assessment for ALS Using Multimodal Dialog Agents: Data Quality, Feasibility and Task Compliance.....	5441
<i>Vanessa Richter, Michael Neumann, Jordan Green, Brian Richburg, Oliver Roesler, Hardik Kothare, Vikram Ramanarayanan</i>	
AudioToken: Adaptation of Text-Conditioned Diffusion Models for Audio-To-Image Generation.....	5446
<i>Guy Yariv, Itai Gat, Lior Wolf, Yossi Adi, Idan Schwartz</i>	
Obstructive Sleep Apnea Screening with Breathing Sounds and Respiratory Effort: A Multimodal Deep Learning Approach.....	5451
<i>Hector E. Romero, Ning Ma, Guy J. Brown, Sam Johnson</i>	
Investigation of Music Emotion Recognition Based on Segmented Semi-Supervised Learning	5456
<i>Yifu Sun, Xulong Zhang, Jianzong Wang, Ning Cheng, Kaiyu Hu, Jing Xiao</i>	

SPEECH SYNTHESIS: MULTILINGUALITY; EVALUATION

The Effects of Input Type and Pronunciation Dictionary Usage in Transfer Learning for Low-Resource Text-To-Speech.....	5461
<i>Phat Do, Matt Coler, Jelske Dijkstra, Esther Klabbers</i>	
Resource-Efficient Fine-Tuning Strategies for Automatic MOS Prediction in Text-To-Speech for Low-Resource Languages	5466
<i>Phat Do, Matt Coler, Jelske Dijkstra, Esther Klabbers</i>	
Robust Feature Decoupling in Voice Conversion by Using Locality-Based Instance Normalization	5471
<i>Yewei Gu, Xianfeng Zhao, Xiaowei Yi</i>	
Zero-Shot Accent Conversion Using Pseudo Siamese Disentanglement Network	5476
<i>Dongya Jia, Qiao Tian, Kainan Peng, Jiaxin Li, Yuanzhe Chen, Mingbo Ma, Yuping Wang, Yuxuan Wang</i>	
Automatic Evaluation of Turn-Taking Cues in Conversational Speech Synthesis	5481
<i>Erik Ekstedt, Siyang Wang, Éva Székely, Joakim Gustafson, Gabriel Skantze</i>	
GenerTTS: Pronunciation Disentanglement for Timbre and Style Generalization in Cross-Lingual Text-To-Speech.....	5486
<i>Yahuan Cong, Haoyu Zhang, Haopeng Lin, Shichao Liu, Chunfeng Wang, Yi Ren, Xiang Yin, Zejun Ma</i>	
Analysis of Mean Opinion Scores in Subjective Evaluation of Synthetic Speech Based on Tail Probabilities.....	5491
<i>Yusuke Yasuda, Tomoki Toda</i>	
LibriTTS-R: A Restored Multi-Speaker Text-To-Speech Corpus.....	5496
<i>Yuma Koizumi, Heiga Zen, Shigeki Karita, Yifan Ding, Kohei Yatabe, Nobuyuki Morioka, Michiel Bacchiani, Yu Zhang, Wei Han, Ankur Bapna</i>	
UniFLG: Unified Facial Landmark Generator from Text Or Speech	5501
<i>Kentaro Mitsui, Yukiya Hono, Kei Sawada</i>	
XPhoneBERT: A Pre-Trained Multilingual Model for Phoneme Representations for Text-To-Speech	5506
<i>Linh The Nguyen, Thinh Pham, Dat Quoc Nguyen</i>	

CIArTTS: An Open-Source Classical Arabic Text-To-Speech Corpus.....	5511
<i>Ajinkya Kulkarni, Atharva Kulkarni, Sara Abedalmon'Em Mohammad Shatnawi, Hanan Aldarmaki</i>	
Diffusion-Based Accent Modelling in Speech Synthesis	5516
<i>Kamil Deja, Georgi Tinchev, Marta Czarnowska, Marius Cotescu, Jasha Droppo</i>	
Multilingual Text-To-Speech Synthesis for Turkic Languages Using Transliteration.....	5521
<i>Rustem Yeshpanov, Saida Mussakhojayeva, Yerbolat Khassanov</i>	
CVTE-Poly: A New Benchmark for Chinese Polyphone Disambiguation.....	5526
<i>Siheng Zhang, Xingjun Tan, Yanqiang Lei, Xianxiang Wang, Zhizhong Zhang, Yuan Xie</i>	
Improving Bilingual TTS Using Language and Phonology Embedding with Embedding Strength Modulator.....	5531
<i>Fengyu Yang, Jian Luan, Meng Meng, Yujun Wang</i>	
High-Quality Automatic Voice Over with Accurate Alignment: Supervision Through Self-Supervised Discrete Speech Units.....	5536
<i>Junchen Lu, Berrak Sisman, Mingyang Zhang, Haizhou Li</i>	
PronScribe: Highly Accurate Multimodal Phonemic Transcription from Speech and Text	5541
<i>Yang Yu, Matthew Perez, Ankur Bapna, Fadi Haik, Siamak Tazari, Yu Zhang</i>	
Expressive Machine Dubbing Through Phrase-Level Cross-Lingual Prosody Transfer	5546
<i>Jakub Swiatkowski, Duo Wang, Mikolaj Babianski, Giuseppe Coccia, Patrick Lumban Tobing, Ravichander Vipperla, Viacheslav Klimkov, Vincent Pollet</i>	
Why We Should Report the Details in Subjective Evaluation of TTS More Rigorously	5551
<i>Cheng-Han Chiang, Wei-Ping Huang, Hung-Yi Lee</i>	
Speaker-Independent Neural Formant Synthesis.....	5556
<i>Pablo Pérez Zarazaga, Zofia Malisz, Gustav Eje Henter, Lauri Juvela</i>	
CALLS: Japanese Empathetic Dialogue Speech Corpus of Complaint Handling and Attentive Listening in Customer Center.....	5561
<i>Yuki Saito, Eiji Imori, Shinnosuke Takamichi, Kentaro Tachibana, Hiroshi Saruwatari</i>	
SASPEECH: A Hebrew Single Speaker Dataset for Text to Speech and Voice Conversion.....	5566
<i>Orian Sharoni, Roe Shenberg, Erica Cooper</i>	

Author Index