

# IEEE/ACM TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING

A PUBLICATION OF THE IEEE SIGNAL PROCESSING SOCIETY



Association for  
Computing Machinery

[www.signalprocessingsociety.org](http://www.signalprocessingsociety.org)

Indexed in PubMed® and MEDLINE®, products of the United States National Library of Medicine



2021

VOLUME 29

ITASFA

(ISSN 2329-9290)

REGULAR PAPERS

Hierarchical Regulated Iterative Network for Joint Task of Music Detection and Music Relative Loudness Estimation . . . . .	<i>B. Jia, J. Lv, X. Peng, Y. Chen, and S. Yang</i>	1
Novel Architectures for Unsupervised Information Bottleneck Based Speaker Diarization of Meetings . . . . .	<i>N. Dawalatabad, S. Madikeri, C. C. Sekhar, and H. A. Murthy</i>	14
Block-Based High Performance CNN Architectures for Frame-Level Overlapping Speech Detection . . . . .	<i>M. Yousefi and J. H. L. Hansen</i>	28
A Deep Adaptation Network for Speech Enhancement: Combining a Relativistic Discriminator With Multi-Kernel Maximum Mean Discrepancy . . . . .	<i>J. Cheng, R. Liang, Z. Liang, L. Zhao, C. Huang, and B. Schuller</i>	41
Comparison of Artificial Neural Network Types for Infant Vocalization Classification . . . . .	<i>F. Anders, M. Hlawitschka, and M. Fuchs</i>	54
Harmonic-Temporal Factor Decomposition for Unsupervised Monaural Separation of Harmonic Sounds . . . . .	<i>T. Nakamura and H. Kameoka</i>	68
Computation of Spherical Harmonic Representations of Source Directivity Based on the Finite-Distance Signature . . . . .	<i>J. Ahrens and S. Bilbao</i>	83
Improving Automatic Speech Recognition and Speech Translation via Word Embedding Prediction . . . . .	<i>S.-P. Chuang, A. H. Liu, T.-W. Sung, and H.-yi Lee</i>	93
A Cross-Entropy-Guided Measure (CEGM) for Assessing Speech Recognition Performance and Optimizing DNN-Based Speech Enhancement . . . . .	<i>L. Chai, J. Du, Q.-F. Liu, and C.-H. Lee</i>	106
Passive Geometry Calibration for Microphone Arrays Based on Distributed Damped Newton Optimization . . . . .	<i>D. Hu, Z. Chen, and F. Yin</i>	118
An Overview of Voice Conversion and Its Challenges: From Statistical Modeling to Deep Learning . . . . .	<i>B. Sisman, J. Yamagishi, S. King, and H. Li</i>	132

(Contents Continued on Page i)



---

Steering Study of Linear Differential Microphone Arrays . . . . .	158
Proportionate Adaptive Filtering Algorithms Derived Using an Iterative Reweighting Framework . . . . .	171
A Novel Approach for Improved Noise Reduction Performance in Feed-Forward Active Noise Control Systems With (Loudspeaker) Saturation Non-Linearity in the Secondary Path . . . . .	187
Gated Recurrent Fusion With Joint Training Framework for Robust End-to-End Speech Recognition . . . . .	198
Speech Intelligibility Prediction Using Spectro-Temporal Modulation Analysis . . . . .	210
On the Design of Sparse Arrays With Frequency-Invariant Beam Pattern . . . . .	226
A Room Compensation Method by Modification of Reverberant Audio Objects . . . . .	239
Multiple Source Direction of Arrival Estimations Using Relative Sound Pressure Based MUSIC . . . . .	253
The Temporal Limits Encoder as a Sound Coding Strategy for Bilateral Cochlear Implants . . . . .	265
Exploiting Morphological and Phonological Features to Improve Prosodic Phrasing for Mongolian Speech Synthesis . . . . .	274
Multiple Circular Arrays of Vector Sensors for Real-Time Sound Field Analysis . . . . .	286
Robust Sound Source Tracking Using SRP-PHAT and 3D Convolutional Neural Networks . . . . .	300
Directly Comparing the Listening Strategies of Humans and Machines . . . . .	312
Auxiliary Networks for Joint Speaker Adaptation and Speaker Change Detection . . . . .	324
Multiple Acoustic Source Localization in Microphone Array Networks . . . . .	334
Tackling Perception Bias in Unsupervised Phoneme Discovery Using DPGMM-RNN Hybrid Model and Functional Load . . . . .	348
Fast Generation of Sound Zones Using Variable Span Trade-Off Filters in the DFT-Domain . . . . .	363
Multi-Source DOA Estimation in Reverberant Environments by Jointing Detection and Modeling of Time-Frequency Points . . . . .	379
Speech Enhancement Based on Modulation-Domain Parametric Multichannel Kalman Filtering . . . . .	393
A Knowledge Graph Embedding Approach for Metaphor Processing . . . . .	406
Estimation Reliability Function Assisted Sound Source Localization With Enhanced Steering Vector Phase Difference . . . . .	421
Room Acoustical Parameter Estimation From Room Impulse Responses Using Deep Neural Networks . . . . .	436
Affine Projection Algorithm Over Acoustic Sensor Networks for Active Noise Control . . . . .	448
Performance Analysis of the Extended Binaural MVDR Beamformer With Partial Noise Estimation . . . . .	462
Ensemble Bag-of-Audio-Words Representation Improves Paralinguistic Classification Accuracy . . . . .	477
Room Impulse Response Reshaping and Crosstalk Cancellation Using Convex Optimization . . . . .	489
Investigating Typed Syntactic Dependencies for Targeted Sentiment Classification Using Graph Attention Neural Network . . . . .	503
Speech Enhancement via Attention Masking Network (SEAMNET): An End-to-End System for Joint Suppression of Noise and Reverberation . . . . .	515
Voice Jitter Estimation Using High-Order Synchronizing Operators . . . . .	527

---

Speaker Separation Using Speaker Inventories and Estimated Speech . . . . .	
. . . . . <i>P. Wang, Z. Chen, D. L. Wang, J. Li, and Y. Gong</i>	537
On the Distribution of Speaker Verification Scores: Generative Models for Unsupervised Calibration . . . . .	<i>S. Cumani</i> 547
Acoustic Measure of Vocal Strain Based on Glottal Airflow Periodicity . . . . .	<i>Y.-R. Chien and J. Guanason</i> 563
RARS: Recognition of Audio Recording Source Based on Residual Neural Network . . . . .	
. . . . . <i>X. Shen, X. Shao, Q. Ge, and L. Liu</i>	575
Learning to Generate Explainable Plots for Neural Story Generation . . . . .	
. . . . . <i>G. Chen, Y. Liu, H. Luan, M. Zhang, Q. Liu, and M. Sun</i>	585
A New Class of Differential Beamformers . . . . .	<i>W. Yang, J. Benesty, G. Huang, and J. Chen</i> 594
Multichannel Blind Source Separation Based on Evanescent-Region-Aware Non-Negative Tensor Factorization in Spherical Harmonic Domain . . . . .	<i>Y. Mitsufuji, N. Takamune, S. Koyama, and H. Saruwatari</i> 607
Robust Constrained MFMVDR Filters for Single-Channel Speech Enhancement Based on Spherical Uncertainty Set . . . . .	
. . . . . <i>D. Fischer and S. Doclo</i>	618
Differential Beamforming From the Beampattern Factorization Perspective . . . . .	
. . . . . <i>X. Zhao, J. Benesty, J. Chen, and G. Huang</i>	632
Preordering Encoding on Transformer for Translation . . . . .	<i>Y. Kawara, C. Chu, and Y. Arase</i> 644
Many-to-Many Voice Transformer Network . . . . .	<i>H. Kameoka, W.-C. Huang, K. Tanaka, T. Kaneko, N. Hojo, and T. Toda</i> 656
A Study on Reference Microphone Selection for Multi-Microphone Speech Enhancement . . . . .	
. . . . . <i>J. Zhang, H. Chen, L.-R. Dai, and R. C. Hendriks</i>	671
Overview and Evaluation of Sound Event Localization and Detection in DCASE 2019 . . . . .	
. . . . . <i>A. Politis, A. Mesaros, S. Adavanne, T. Heittola, and T. Virtanen</i>	684
Listening Enhancement in Noisy Environments: Solutions in Time and Frequency Domain . . . . .	
. . . . . <i>M. Niermann and P. Vary</i>	699
Gated Recurrent Context: Softmax-Free Attention for Online Encoder-Decoder Speech Recognition . . . . .	
. . . . . <i>H. Lee, W. H. Kang, S. J. Cheon, H. Kim, and N. S. Kim</i>	710
On Improved Training of CNN for Acoustic Source Localisation . . . . .	<i>E. Vargas, J. R. Hopgood, K. Brown, and K. Subr</i> 720
Deep Normalization for Speaker Vectors . . . . .	<i>Y. Cai, L. Li, A. Abel, X. Zhu, and D. Wang</i> 733
Pretraining Techniques for Sequence-to-Sequence Voice Conversion . . . . .	
. . . . . <i>W.-C. Huang, T. Hayashi, Y.-C. Wu, H. Kameoka, and T. Toda</i>	745
Temporal Dynamics of Workplace Acoustic Scenes: Egocentric Analysis and Prediction . . . . .	
. . . . . <i>A. Jati, A. Nadarajan, R. Peri, K. Mundnich, T. Feng, B. Girault, and S. Narayanan</i>	756
Modeling Future Cost for Neural Machine Translation . . . . .	
. . . . . <i>C. Duan, K. Chen, R. Wang, M. Utiyama, E. Sumita, C. Zhu, and T. Zhao</i>	770
Adaptive Convolution for Semantic Role Labeling . . . . .	<i>K. Munir, H. Zhao, and Z. Li</i> 782
Quasi-Periodic Parallel WaveGAN: A Non-Autoregressive Raw Waveform Generative Model With Pitch-Dependent Dilated Convolution Neural Network . . . . .	<i>Y.-C. Wu, T. Hayashi, T. Okamoto, H. Kawai, and T. Toda</i> 792
Evolving Multi-Resolution Pooling CNN for Monaural Singing Voice Separation . . . . .	
. . . . . <i>W. Yuan, B. Dong, S. Wang, M. Unoki, and W. Wang</i>	807
Generation of Personal Sound Zones With Physical Meaningful Constraints and Conjugate Gradient Method . . . . .	
. . . . . <i>L. Shi, T. Lee, L. Zhang, J. K. Nielsen, and M. G. Christensen</i>	823
On the Robustness of the Superdirective Beamformer . . . . .	<i>X. Chen, J. Benesty, G. Huang, and J. Chen</i> 838
Generating Images From Spoken Descriptions . . . . .	<i>X. Wang, T. Qiao, J. Zhu, A. Hanjalic, and O. Scharenborg</i> 850
Domain-Aware Dialogue State Tracker for Multi-Domain Dialogue Systems . . . . .	<i>V. Balaraman and B. Magnini</i> 866
Exemplar-Based Emotive Speech Synthesis . . . . .	<i>X. Wu, Y. Cao, H. Lu, S. Liu, S. Kang, Z. Wu, X. Liu, and H. Meng</i> 874
Towards Duration Robust Weakly Supervised Sound Event Detection . . . . .	<i>H. Dinkel, M. Wu, and K. Yu</i> 887
Binaural Reproduction Based on Bilateral Ambisonics and Ear-Aligned HRTFs . . . . .	
. . . . . <i>Z. Ben-Hur, D. L. Alon, R. Mehra, and B. Rafaely</i>	901
Synthesis and Analysis-By-Synthesis of Modulated Diplophonic Glottal Area Waveforms . . . . .	
. . . . . <i>P. Aichinger and F. Pernkopf</i>	914
Analysis and Calibration of Lombard Effect and Whisper for Speaker Recognition . . . . .	<i>F. Kelly and J. H.L. Hansen</i> 927
Lyric or Dramatic - Vibrato Analysis for Voice Type Classification in Professional Opera Singers . . . . .	
. . . . . <i>M. Müller, T. Schulz, T. Ermakova, and P. P. Caffier</i>	943
Incorporating Wireless Communication Parameters Into the E-Model Algorithm . . . . .	
. . . . . <i>D. Z. Rodríguez, D. Carrillo, M. A. Ramírez, P. H. J. Nardelli, and S. Möller</i>	956

---

Non-Linear-Echo Based Anti-Collusion Mechanism for Audio Signals . . . . .	969
. . . . . <i>T. Zong, Y. Xiang, I. Natgunanathan, L. Gao, G. Hua, and W. Zhou</i>	
CTNet: Conversational Transformer Network for Emotion Recognition . . . . .	985
. . . . . <i>Z. Lian, B. Liu, and J. Tao</i>	
Neural Machine Translation With Explicit Phrase Alignment . . . . .	1001
. . . . . <i>J. Zhang, H. Luan, M. Sun, F. Zhai, J. Xu, and Y. Liu</i>	
Cognitive Load Estimation From Speech Commands to Simulated Aircraft . . . . .	1011
. . . . . <i>M. Vukovic, M. Stolar, and M. Lech</i>	
Geometry Calibration for Acoustic Transceiver Networks Based on Network Newton Distributed Optimization . . . . .	1023
. . . . . <i>D. Hu, Z. Chen, and F. Yin</i>	
Perceptual-Similarity-Aware Deep Speaker Representation Learning for Multi-Speaker Generative Modeling . . . . .	1033
. . . . . <i>Y. Saito, S. Takamichi, and H. Saruwatari</i>	
Vocal Tract Length Estimation Using Accumulated Means of Formants and Its Effects on Speaker-Normalization . . . . .	1049
. . . . . <i>T. Sakata, N. Ikeda, Y. Ueda, and A. Watanabe</i>	
Modified Magnitude-Phase Spectrum Information for Spoofing Detection . . . . .	1065
. . . . . <i>J. Yang, H. Wang, R. K. Das, and Y. Qian</i>	
Audio-Visual Deep Neural Network for Robust Person Verification . . . . .	1079
. . . . . <i>Y. Qian, Z. Chen, and S. Wang</i>	
Deep Selective Memory Network With Selective Attention and Inter-Aspect Modeling for Aspect Level Sentiment Classification . . . . .	1093
. . . . . <i>P. Lin, M. Yang, and J. Lai</i>	
Improved Acoustic Word Embeddings for Zero-Resource Languages Using Multilingual Transfer . . . . .	1107
. . . . . <i>H. Kamper, Y. Matuskevych, and S. Goldwater</i>	
Audio-Based Piano Performance Evaluation for Beginners With Convolutional Neural Network and Attention Mechanism . . . . .	1119
. . . . . <i>W. Wang, J. Pan, H. Yi, Z. Song, and M. Li</i>	
Quasi-Periodic WaveNet: An Autoregressive Raw Waveform Generative Model With Pitch-Dependent Dilated Convolution Neural Network . . . . .	1134
. . . . . <i>Y.-C. Wu, T. Hayashi, P. L. Tobing, K. Kobayashi, and T. Toda</i>	
Late-Reverberation Synthesis Using Interleaved Velvet-Noise Sequences . . . . .	1149
. . . . . <i>V. Välimäki and K. Prawda</i>	
Multi-Turn Dialogue Reading Comprehension With Pivot Turns and Knowledge . . . . .	1161
. . . . . <i>Z. Zhang, J. Li, and H. Zhao</i>	
Sparsity-Based Audio Declipping Methods: Selected Overview, New Algorithms, and Large-Scale Evaluation . . . . .	1174
. . . . . <i>C. Gaultier, S. Kitić, R. Gribonval, and N. Bertin</i>	
Mixed Source Sound Field Translation for Virtual Binaural Application With Perceptual Validation . . . . .	1188
. . . . . <i>L. Birnie, T. Abhayapala, V. Tourbabin, and P. Samarasinghe</i>	
Meta-Learning With Latent Space Clustering in Generative Adversarial Network for Speaker Diarization . . . . .	1204
. . . . . <i>M. Pal, M. Kumar, R. Peri, T. J. Park, S. H. Kim, C. Lord, S. Bishop, and S. Narayanan</i>	
Sensor Selection for Relative Acoustic Transfer Function Steered Linearly-Constrained Beamformers . . . . .	1220
. . . . . <i>J. Zhang, J. Du, and L.-R. Dai</i>	
Zero-Shot Audio Classification Via Semantic Embeddings . . . . .	1233
. . . . . <i>H. Xie and T. Virtanen</i>	
Phoneme-Unit-Specific Time-Delay Neural Network for Speaker Verification . . . . .	1243
. . . . . <i>X. Chen and C. Bao</i>	
Optimal Output-Constrained Active Noise Control Based on Inverse Adaptive Modeling Leak Factor Estimate . . . . .	1256
. . . . . <i>D. Shi, W.-S. Gan, B. Lam, S. Wen, and X. Shen</i>	
Dense CNN With Self-Attention for Time-Domain Speech Enhancement . . . . .	1270
. . . . . <i>A. Pandey and D. L. Wang</i>	
Knowing Where to Leverage: Context-Aware Graph Convolutional Network With an Adaptive Fusion Layer for Contextual Spoken Language Understanding . . . . .	1280
. . . . . <i>L. Qin, W. Che, M. Ni, Y. Li, and T. Liu</i>	
Transfer Learning From Speech Synthesis to Voice Conversion With Non-Parallel Training Data . . . . .	1290
. . . . . <i>M. Zhang, Y. Zhou, L. Zhao, and H. Li</i>	
Neural Network Adaptation and Data Augmentation for Multi-Speaker Direction-of-Arrival Estimation . . . . .	1303
. . . . . <i>W. He, P. Motlicek, and J.-M. Odobez</i>	
Improving Skip-Gram Embeddings Using BERT . . . . .	1318
. . . . . <i>Y. Wang, L. Cui, and Y. Zhang</i>	
Deep Graph-Based Character-Level Chinese Dependency Parsing . . . . .	1329
. . . . . <i>L. Wu and M. Zhang</i>	
Integrating Knowledge Into End-to-End Speech Recognition From External Text-Only Data . . . . .	1340
. . . . . <i>Y. Bai, J. Yi, J. Tao, Z. Wen, Z. Tian, and S. Zhang</i>	
Convolutional Maximum-Likelihood Distortionless Response Beamforming With Steering Vector Estimation for Robust Speech Recognition . . . . .	1352
. . . . . <i>B. J. Cho and H.-M. Park</i>	
An Overview of Deep-Learning-Based Audio-Visual Speech Enhancement and Separation . . . . .	1368
. . . . . <i>D. Michelsanti, Z.-H. Tan, S.-X. Zhang, Y. Xu, M. Yu, D. Yu, and J. Jensen</i>	
On the Design of Differential Kronecker Product Beamformers . . . . .	1397
. . . . . <i>G. Itzhak, J. Benesty, and I. Cohen</i>	
Partially Matching Projection Decoding Method Evaluation Under Different Playback Conditions . . . . .	1411
. . . . . <i>Z. Ge, L. Li, and T. Qu</i>	
Analyzing Multimodal Sentiment Via Acoustic- and Visual-LSTM With Channel-Aware Temporal Convolution Network . . . . .	1424
. . . . . <i>S. Mai, S. Xing, and H. Hu</i>	



---

A Joint Model for Named Entity Recognition With Sentence-Level Entity Type Attentions . . . . .	1438
. . . . . <i>T. Qian, M. Zhang, Y. Lou, and D. Hua</i>	
Ambisonic Signal Processing DNNs Guaranteeing Rotation, Scale and Time Translation Equivariance . . . . .	1449
. . . . . <i>R. Sato, K. Niwa, and K. Kobayashi</i>	
Iterative Echo Labeling Algorithm With Convex Hull Expansion for Room Geometry Estimation . . . . .	1463
. . . . . <i>S. Park and J.-W. Choi</i>	
Overlapping Speaker Segmentation Using Multiple Hypothesis Tracking of Fundamental Frequency . . . . .	1479
. . . . . <i>A. O. T. Hogg, C. Evers, A. H. Moore, and P. A. Naylor</i>	
Controlling Elevation and Azimuth Beamwidths With Concentric Circular Microphone Arrays . . . . .	1491
. . . . . <i>R. Sharma, I. Cohen, and B. Berdugo</i>	
A Multiple-Integration Encoder for Multi-Turn Text-to-SQL Semantic Parsing . . . . .	1503
. . . . . <i>R.-Z. Wang, Z.-H. Ling, J.-B. Zhou, and Y. Hu</i>	
Bayesian Learning of LF-MMI Trained Time Delay Neural Networks for Speech Recognition . . . . .	1514
. . . . . <i>S. Hu, X. Xie, S. Liu, J. Yu, Z. Ye, M. Geng, X. Liu, and H. Meng</i>	
Objective Measures of Perceptual Audio Quality Reviewed: An Evaluation of Their Application Domain Dependence . . . . .	1530
. . . . . <i>M. Torcoli, T. Kastner, and J. Herre</i>	
Voice Activity Detection in the Wild: A Data-Driven Approach Using Teacher-Student Training . . . . .	1542
. . . . . <i>H. Dinkel, S. Wang, X. Xu, M. Wu, and K. Yu</i>	
Detection of Multiple Steganography Methods in Compressed Speech Based on Code Element Embedding, Bi-LSTM and CNN With Attention Mechanisms . . . . .	1556
. . . . . <i>S. Li, J. Wang, P. Liu, M. Wei, and Q. Yan</i>	
Deformable Self-Attention for Text Classification . . . . .	1570
. . . . . <i>Q. Ma, J. Yan, Z. Lin, L. Yu, and Z. Chen</i>	
Extracting and Predicting Word-Level Style Variations for Speech Synthesis . . . . .	1582
. . . . . <i>Y.-J. Zhang and Z.-H. Ling</i>	
Exploiting Temporal Context in CNN Based Multisource DOA Estimation . . . . .	1594
. . . . . <i>A. Bohlender, A. Spriet, W. Tirry, and N. Madhu</i>	
Determined BSS Based on Time-Frequency Masking and Its Application to Harmonic Vector Analysis . . . . .	1609
. . . . . <i>K. Yatabe and D. Kitamura</i>	
TutorNet: Towards Flexible Knowledge Distillation for End-to-End Speech Recognition . . . . .	1626
. . . . . <i>J. W. Yoon, H. Lee, H. Y. Kim, W. I. Cho, and N. S. Kim</i>	
Self-Supervised Representation Learning With Path Integral Clustering for Speaker Diarization . . . . .	1639
. . . . . <i>P. Singh and S. Ganapathy</i>	
A Graph-to-Sequence Learning Framework for Summarizing Opinionated Texts . . . . .	1650
. . . . . <i>P. Wei, J. Zhao, and W. Mao</i>	
Near-Field Superdirectivity: An Analytical Perspective . . . . .	1661
. . . . . <i>D. Y. Levin, S. Markovich-Golan, and S. Gannot</i>	
Speech Emotion Recognition Considering Nonverbal Vocalization in Affective Conversations . . . . .	1675
. . . . . <i>J.-H. Hsu, M.-H. Su, C.-H. Wu, and Y.-H. Chen</i>	
Time-Domain Audio Source Separation With Neural Networks Based on Multiresolution Analysis . . . . .	1687
. . . . . <i>T. Nakamura, S. Kozuka, and H. Saruwatari</i>	
FSPRM: A Feature Subsequence Based Probability Representation Model for Chinese Word Embedding . . . . .	1702
. . . . . <i>Y. Zhang, Y. Liu, J. Zhu, and X. Wu</i>	
Any-to-Many Voice Conversion With Location-Relative Sequence-to-Sequence Modeling . . . . .	1717
. . . . . <i>S. Liu, Y. Cao, D. Wang, X. Wu, X. Liu, and H. Meng</i>	
An Optimal Envelope-Based Noise Reduction Method for Cochlear Implants: An Upper Bound Performance Investigation . . . . .	1729
. . . . . <i>R. A. Chiea, M. H. Costa, and J. A. Cordioli</i>	
Adaptive Adapters: An Efficient Way to Incorporate BERT Into Neural Machine Translation . . . . .	1740
. . . . . <i>J. Guo, Z. Zhang, L. Xu, B. Chen, and E. Chen</i>	
Group Communication With Context Codec for Lightweight Source Separation . . . . .	1752
. . . . . <i>Y. Luo, C. Han, and N. Mesgarani</i>	
Hierarchical Neighbor Propagation With Bidirectional Graph Attention Network for Relation Prediction . . . . .	1762
. . . . . <i>Z. Xie, R. Zhu, J. Liu, G. Zhou, and J. X. Huang</i>	
Beamforming with Cube Microphone Arrays Via Kronecker Product Decompositions . . . . .	1774
. . . . . <i>X. Wang, J. Benesty, J. Chen, G. Huang, and I. Cohen</i>	
Towards Model Compression for Deep Learning Based Speech Enhancement . . . . .	1785
. . . . . <i>K. Tan and D. L. Wang</i>	
Nonlinear Spatial Filtering in Multichannel Speech Enhancement . . . . .	1795
. . . . . <i>K. Tesch and T. Gerkmann</i>	
Expressive TTS Training With Frame and Style Reconstruction Loss . . . . .	1806
. . . . . <i>R. Liu, B. Sisman, G. Gao, and H. Li</i>	
Chinese Lexical Simplification . . . . .	1819
. . . . . <i>J. Qiang, X. Lu, Y. Li, Y. Yuan, and X. W.</i>	
Two Heads are Better Than One: A Two-Stage Complex Spectral Mapping Approach for Monaural Speech Enhancement . . . . .	1829
. . . . . <i>A. Li, W. Liu, C. Zheng, C. Fan, and X. Li</i>	

---

Weighted Orthogonal Vector Rejection Method for Loudspeaker-Based Binaural Audio Reproduction . . . . .	<i>E. C. Hamdan and F. M. Fazi</i>	1844
Deep Learning Based Real-Time Speech Enhancement for Dual-Microphone Mobile Phones . . . . .	<i>K. Tan, X. Zhang, and D. L. Wang</i>	1853
Indoor Multi-Speaker Localization Based on Bayesian Nonparametrics in the Circular Harmonic Domain . . . . .	<i>K. SongGong, H. Chen, and W. Wang</i>	1864
Double-Cross-Correlation Processing for Blind Sampling-Rate and Time-Offset Estimation . . . . .	<i>A. Chinaev, P. Thüine, and G.ENZNER</i>	1881
Fast End-to-End Speech Recognition Via Non-Autoregressive Models and Cross-Modal Knowledge Transferring From BERT . . . . .	<i>Y. Bai, J. Yi, J. Tao, Z. Tian, Z. Wen, and S. Zhang</i>	1897
Multimodal Representations for Synchronized Speech and Real-Time MRI Video Processing . . . . .	<i>Ö.ykü. D. Köse and M. Saraçlar</i>	1912
The Detection of Parkinson’s Disease From Speech Using Voice Source Information . . . . .	<i>N. P. Narendra, B. Schuller, and P. Alku</i>	1925
SNR-Based Features and Diverse Training Data for Robust DNN-Based Speech Enhancement . . . . .	<i>R. Rehr and T. Gerkmann</i>	1937
A Joint Diagonalization Based Efficient Approach to Underdetermined Blind Audio Source Separation Using the Multichannel Wiener Filter . . . . .	<i>N. Ito, R. Ikeshita, H. Sawada, and T. Nakatani</i>	1950
Second-Order Semantic Role Labeling With Global Structural Refinement . . . . .	<i>H. Fei, S. Wu, Y. Ren, and D. Ji</i>	1966
F0 Perturbation Due to Articulatory Movements: Filtering, Characterization and Applications . . . . .	<i>H. M. Torres, M. Güemes, J. A. Gurlekian, and D. A. Evin</i>	1977
Receptive Field Regularization Techniques for Audio Classification and Tagging With Deep Convolutional Neural Networks . . . . .	<i>K. Koutini, H. Eghbal-zadeh, and G. Widmer</i>	1987
Multi-microphone Complex Spectral Mapping for Utterance-wise and Continuous Speech Separation . . . . .	<i>Z.-Q. Wang, P. Wang, and D. L. Wang</i>	2001
Relation Extraction in Dialogues: A Deep Learning Model Based on the Generality and Specialty of Dialogue Text . . . . .	<i>M. Zhou, D. Ji, and F. Li</i>	2015
Domain-Shift Conditioning Using Adaptable Filtering Via Hierarchical Embeddings for Robust Chinese Spell Check . . . . .	<i>M. Nguyen, G. H. Ngo, and N. F. Chen</i>	2027
The Effect of Partial Time-Frequency Masking of the Direct Sound on the Perception of Reverberant Speech . . . . .	<i>L. Madmoni, S. Tibor, I. Nelken, and B. Rafaely</i>	2037
Corpus-Aware Graph Aggregation Network for Sequence Labeling . . . . .	<i>H. Chen, Q. Ma, L. Yu, Z. Lin, and J. Yan</i>	2048
Towards Robust Speech Super-Resolution . . . . .	<i>H. Wang and D. L. Wang</i>	2058
Audio-Visual Multi-Channel Integration and Recognition of Overlapped Speech . . . . .	<i>J. Yu, S.-X. Zhang, B. Wu, S. Liu, S. Hu, M. Geng, X. Liu, H. Meng, and D. Yu</i>	2067
Conditioned Source Separation for Musical Instrument Performances . . . . .	<i>O. Slizovskaia, G. Haro, and E. Gómez</i>	2083
Bayesian Learning for Deep Neural Network Adaptation . . . . .	<i>X. Xie, X. Liu, T. Lee, and L. Wang</i>	2096
Nearest Kronecker Product Decomposition Based Linear-in-The-Parameters Nonlinear Filters . . . . .	<i>S. S. Bhattacharjee and N. V. George</i>	2111
A Multi-Agent Communication Based Model for Nested Named Entity Recognition . . . . .	<i>C. Li, G. Wang, J. Cao, and Y. Cai</i>	2123
Blind Separation for Multiple Moving Sources With Labeled Random Finite Sets . . . . .	<i>J. Ong, B. T. Vo, and S. Nordholm</i>	2137
PROTOTYPE-TO-STYLE: Dialogue Generation With Style-Aware Editing on Retrieval Memory . . . . .	<i>Y. Su, Y. Wang, D. Cai, S. Baker, A. Korhonen, and N. Collier</i>	2152
A Wave Digital Newton-Raphson Method for Virtual Analog Modeling of Audio Circuits with Multiple One-Port Nonlinearities . . . . .	<i>A. Bernardini, E. Bozzo, F. Fontana, and A. Sarti</i>	2162
Proximal Normalized Subband Adaptive Filtering for Acoustic Echo Cancellation . . . . .	<i>G. Guo, Y. Yu, R. C. de Lamare, Z. Zheng, L. Lu, and Q. Cai</i>	2174
Audibility of Group-Delay Equalization . . . . .	<i>J. Liski, A. Mäkitvirta, and V. Välimäki</i>	2189
Corruption Is Not All Bad: Incorporating Discourse Structure Into Pre-Training via Corruption for Essay Scoring . . . . .	<i>F. S. Mim, N. Inoue, P. Reisert, H. Ouchi, and K. Inui</i>	2202
Graph-Based Clustering of Dolphin Whistles . . . . .	<i>D. Kipnis and R. Diamant</i>	2216
Language-Independent Approach for Automatic Computation of Vowel Articulation Features in Dysarthric Speech Assessment . . . . .	<i>Y. Liu, N. Penttilä, T. Ihalainen, J. Lintula, R. Convey, and O. Räsänen</i>	2228
Impulsive Noise Detection for Speech Enhancement in HHT Domain . . . . .	<i>C. Medina, R. Coelho, and L. Zăo</i>	2244

---

A Novel Loss Function and Training Strategy for Noise-Robust Keyword Spotting . . . . .	2254
Recent Progress in the CUHK Dysarthric Speech Recognition System . . . . .	2267
Desynchronization Attacks Resilient Watermarking Method Based on Frequency Singular Value Coefficient Modification . . . . .	2282
Sparse Representations With Legendre Kernels for DOA Estimation and Acoustic Source Separation . . . . .	2296
DNN-Based Mask Estimation for Distributed Speech Enhancement in Spatially Unconstrained Microphone Arrays . . . . .	2310
Hybrid Speech and Text Analysis Methods for Speaker Change Detection . . . . .	2324
Multi-Task Sequence Tagging for Emotion-Cause Pair Extraction Via Tag Distribution Refinement . . . . .	2339
TERA: Self-Supervised Learning of Transformer Encoder Representation for Speech . . . . .	2351
Converting Foreign Accent Speech Without a Reference . . . . .	2367
Phoneme Level Lyrics Alignment and Text-Informed Singing Voice Separation . . . . .	2382
High-Order Pair-Wise Aspect and Opinion Terms Extraction With Edge-Enhanced Syntactic Graph Convolution . . . . .	2396
Sentiment Time Series Calibration for Event Detection . . . . .	2407
Learning Context-Aware Convolutional Filters for Implicit Discourse Relation Classification . . . . .	2421
Editorial: Special Issue on the Eighth Dialog System Technology Challenge . . . . .	2434
Randomly Wired Network Based on RoBERTa and Dialog History Attention for Response Selection . . . . .	2437
Deep Contextualized Utterance Representations for Response Selection and Dialogue Analysis . . . . .	2443
End-to-End Recurrent Cross-Modality Attention for Video Dialogue . . . . .	2456
Conversational Semantic Role Labeling . . . . .	2465
Bridging Text and Video: A Universal Multimodal Transformer for Audio-Visual Scene-Aware Dialog . . . . .	2476
GRT: Generative-Retrieval Transformers for Data-Efficient Dialogue Domain Adaptation . . . . .	2484
Domain Adaptive Meta-Learning for Dialogue State Tracking . . . . .	2493
D-Score: Holistic Dialogue Evaluation Without Reference . . . . .	2502
HDRS: Hindi Dialogue Restaurant Search Corpus for Dialogue State Tracking in Task-Oriented Environment . . . . .	2517
Overview of the Eighth Dialog System Technology Challenge: DSTC8 . . . . .	2529
Label and Context Augmentation for Response Selection at DSTC8 . . . . .	2541
History Reuse and Bag-of-Words Loss for Long Summary Generation . . . . .	2551
PhaseDCN: A Phase-Enhanced Dual-Path Dilated Convolutional Network for Single-Channel Speech Enhancement . . . . .	2561
Guided Generative Adversarial Neural Network for Representation Learning and Audio Generation Using Fewer Labelled Audio Data . . . . .	2575
Gamma Boltzmann Machine for Audio Modeling . . . . .	2591
Attending From Foresight: A Novel Attention Mechanism for Neural Machine Translation . . . . .	2606
Information Fusion in Attention Networks Using Adaptive and Multi-Level Factorized Bilinear Pooling for Audio-Visual Emotion Recognition . . . . .	2617
Learning Cross-Lingual Mappings in Imperfectly Isomorphic Embedding Spaces . . . . .	2630

---

UnitNet: A Sequence-to-Sequence Acoustic Model for Concatenative Speech Synthesis . . . . .	<i>X. Zhou, Z.-H. Ling, and L.-R. Dai</i>	2643
Multi-Tone Phase Coding of Interaural Time Difference for Sound Source Localization With Spiking Neural Networks . . . . .	<i>Z. Pan, M. Zhang, J. Wu, J. Wang, and H. Li</i>	2656
FifthNet: Structured Compact Neural Networks for Automatic Chord Recognition . . . . .	<i>K. O'Hanlon and M. B. Sandler</i>	2671
Estimation of Spectral Notches From Pinna Meshes: Insights From a Simple Computational Model . . . . .	<i>S. Spagnol, R. Miccini, M. G. Onofrei, R. Unnthorsson, and S. Serafin</i>	2683
Target Speaker Verification With Selective Auditory Attention for Single and Multi-Talker Speech . . . . .	<i>C. Xu, W. Rao, J. Wu, and H. Li</i>	2696
Minimum Processing Beamforming . . . . .	<i>A. Zahedi, M. S. Pedersen, J. Østergaard, T. U. Christiansen, L. Bramsløw, and J. Jensen</i>	2710
Multichannel Iterative Noise Reduction Filters in the Short-Time-Fourier-Transform Domain Based on Kronecker Product Decomposition . . . . .	<i>X. Wang, J. Chen, X. Chen, J. Guo, and Q. Xiang</i>	2725
Robust $Q$ -Gradient Subband Adaptive Filter for Nonlinear Active Noise Control . . . . .	<i>K.-L. Yin, Y.-F. Pu, and L. Lu</i>	2741
Monaural Speech Separation Using Speaker Embedding From Preliminary Separation . . . . .	<i>J. Byun and J. W. Shin</i>	2753
On the Design of 3D Steerable Beamformers With Uniform Concentric Circular Microphone Arrays . . . . .	<i>X. Zhao, G. Huang, J. Chen, and J. Benesty</i>	2764
A Unified Target-Oriented Sequence-to-Sequence Model for Emotion-Cause Pair Extraction . . . . .	<i>Z. Cheng, Z. Jiang, Y. Yin, N. Li, and Q. Gu</i>	2779
Robust Voice Feature Selection Using Interval Type-2 Fuzzy AHP for Automated Diagnosis of Parkinson's Disease . . . . .	<i>H. Azadi, M.-R. Akbarzadeh-T, H.-R. Kobravi, and A. Shoeibi</i>	2792
Sinsy: A Deep Neural Network-Based Singing Voice Synthesis System . . . . .	<i>Y. Hono, K. Hashimoto, K. Oura, Y. Nankaku, and K. Tokuda</i>	2803
Multi-Granularity Sequence Alignment Mapping for Encoder-Decoder Based End-to-End ASR . . . . .	<i>J. Tang, J. Zhang, Y. Song, I. McLoughlin, and L.-R. Dai</i>	2816
Exploiting Translation Model for Parallel Corpus Mining . . . . .	<i>C. Leong, X. Liu, D. F. Wong, and L. S. Chao</i>	2829
Wavesplit: End-to-End Speech Separation by Speaker Clustering . . . . .	<i>N. Zeghidour and D. Grangier</i>	2840
Learning Waveform-Based Acoustic Models Using Deep Variational Convolutional Neural Networks . . . . .	<i>D. Oglic, Z. Cvetkovic, and P. Sollich</i>	2850
Privacy-Preserving Audio Classification Using Variational Information Feature Extraction . . . . .	<i>A. Nelus and R. Martin</i>	2864
Recurrent Neural Networks and Acoustic Features for Frame-Level Signal-to-Noise Ratio Estimation . . . . .	<i>H. Li, D. L. Wang, X. Zhang, and G. Gao</i>	2878
Generating Responses With a Given Syntactic Pattern in Chinese Dialogues . . . . .	<i>Y. Zhou, X. Zheng, and X. Huang</i>	2888
Binaural Auralization of Microphone Array Room Impulse Responses Using Causal Wiener Filtering . . . . .	<i>V. Gunnarsson and M. Sternad</i>	2899
Sensor Imperfection Tolerance Analysis of Robust Linear Differential Microphone Arrays . . . . .	<i>Z. Chen, H. Chen, and Q. Tu</i>	2915
CSS-LM: A Contrastive Framework for Semi-Supervised Fine-Tuning of Pre-Trained Language Models . . . . .	<i>Y. Su, X. Han, Y. Lin, Z. Zhang, Z. Liu, P. Li, J. Zhou, and M. Sun</i>	2930
A Causality-Constrained Frequency-Domain Least-Squares Filter Design Method for Crosstalk Cancellation . . . . .	<i>T. Kabzinski and P. Jax</i>	2942
CTC-Based Learning of Chroma Features for Score-Audio Music Retrieval . . . . .	<i>F. Zalkow and M. Müller</i>	2957
Multi-Branch Convolutional Macaron net for Sound Event Detection . . . . .	<i>T. K. Chan and C. S. Chin</i>	2972
FluentNet: End-to-End Detection of Stuttered Speech Disfluencies With Deep Learning . . . . .	<i>T. Kourkounakis, A. Hajavi, and A. Etemad</i>	2986
Multi-Metric Optimization Using Generative Adversarial Networks for Near-End Speech Intelligibility Enhancement . . . . .	<i>H. Li and J. Yamagishi</i>	3000
Predict-Then-Decide: A Predictive Approach for Wait or Answer Task in Dialogue Systems . . . . .	<i>Z. Lin, S. Cui, G. Li, X. Kang, F. Ji, F. Li, Z. Zhao, H. Chen, and Y. Zhang</i>	3012
Localization Based on Enhanced Low Frequency Interaural Level Difference . . . . .	<i>M. Calis, S. van de Par, R. Heusdens, and R. C. Hendriks</i>	3025
Native-Nonnative Voice Conversion by Residual Warping in a Sparse, Anchor-Based Representation . . . . .	<i>C. Liberatore</i>	3040
Spatial Active Noise Control Based on Kernel Interpolation of Sound Field . . . . .	<i>S. Koyama, J. Brunnström, H. Ito, N. Ueno, and H. Saruwatari</i>	3052
LSBert: Lexical Simplification Based on BERT . . . . .	<i>J. Qiang, Y. Li, Y. Zhu, Y. Yuan, Y. Shi, and X. Wu</i>	3064



---

Contrastive Information Extraction With Generative Transformer . . . . .	3077
. . . . . <i>N. Zhang, H. Ye, S. Deng, C. Tan, M. Chen, S. Huang, F. Huang, and H. Chen</i>	
Minimum-Volume Multichannel Nonnegative Matrix Factorization for Blind Audio Source Separation . . . . .	3089
. . . . . <i>J. Wang, S. Guan, S. Liu, and X.-L. Zhang</i>	
A Room Impulse Response Measurement Method Robust Towards Nonlinearities Based on Orthogonal Periodic Sequences . . . . .	3104
. . . . . <i>A. Carini, S. Cecchi, A. Terenzi, and S. Orcioni</i>	
Quantization-Aware Binaural MWF Based Noise Reduction Incorporating External Wireless Devices . . . . .	3118
. . . . . <i>J. Zhang and C. Li</i>	
Knowledge Enhanced Fact Checking and Verification . . . . .	3132
. . . . . <i>B. Zhu, X. Zhang, M. Gu, and Y. Deng</i>	
A Superfast Toeplitz Matrix Inversion Method for Single- and Multi-Channel Inverse Filters and Its Application to Room Equalization . . . . .	3144
. . . . . <i>M. A. Poletti and P. D. Teal</i>	
Detecting Source Contextual Barriers for Understanding Neural Machine Translation . . . . .	3158
. . . . . <i>G. Li, L. Liu, C. Zhu, R. Wang, T. Zhao, and S. Shi</i>	
Audio-Aware Spoken Multiple-Choice Question Answering With Pre-Trained Language Models . . . . .	3170
. . . . . <i>C.-C. Kuo, K.-Y. Chen, and S.-B. Luo</i>	
Addressing Extraction and Generation Separately: Keyphrase Prediction With Pre-Trained Language Models . . . . .	3180
. . . . . <i>R. Liu, Z. Lin, and W. Wang</i>	
Sarcasm Detection with Commonsense Knowledge . . . . .	3192
. . . . . <i>J. Li, H. Pan, Z. Lin, P. Fu, and W. Wang</i>	
Keyword Search Using Attention-Based End-to-End ASR and Frame-Synchronous Phoneme Alignments . . . . .	3202
. . . . . <i>R. Yang, G. Cheng, H. Miao, T. Li, P. Zhang, and Y. Yan</i>	
Flexibly Focusing on Supporting Facts, Using Bridge Links, and Jointly Training Specialized Modules for Multi-Hop Question Answering . . . . .	3216
. . . . . <i>T. Alkhaldi, C. Chu, and S. Kurohashi</i>	
An Efficient Filter Bank Structure for Adaptive Notch Filtering and Applications . . . . .	3226
. . . . . <i>W. Wu, Y. Xiao, J. Lin, L. Ma, and K. Khorasani</i>	
Synthesizing Spoken Descriptions of Images . . . . .	3242
. . . . . <i>X. Wang, J. van der Hout, J. Zhu, M. Hasegawa-Johnson, and O. Scharenborg</i>	
Enhancement of Noisy Reverberant Speech Using Polynomial Matrix Eigenvalue Decomposition . . . . .	3255
. . . . . <i>V. W. Neo, C. Evers, and P. A. Naylor</i>	
Wave Digital Modeling and Implementation of Nonlinear Audio Circuits With Nullors . . . . .	3267
. . . . . <i>R. Giampiccolo, M. G. de Bari, A. Bernardini, and A. Sarti</i>	
Speech Emotion Recognition Using Sequential Capsule Networks . . . . .	3280
. . . . . <i>X. Wu, Y. Cao, H. Lu, S. Liu, D. Wang, Z. Wu, X. Liu, and H. Meng</i>	
PSLA: Improving Audio Tagging With Pretraining, Sampling, Labeling, and Aggregation . . . . .	3292
. . . . . <i>Y. Gong, Y.-A. Chung, and J. Glass</i>	
Review and Arrange: Curriculum Learning for Natural Language Understanding . . . . .	3307
. . . . . <i>L. Zhang, Z. Mao, B. Xu, Q. Wang, and Y. Zhang</i>	
Automatic Detection of Affective Flattening in Schizophrenia: Acoustic Correlates to Sound Waves and Auditory Perception . . . . .	3321
. . . . . <i>F. He, L. He, J. Zhang, Y. Li, and X. Xiong</i>	
Efficient Learning Approach for Pronominal Anaphora and Ellipsis Identification and Resolution in Arabic Texts . . . . .	3335
. . . . . <i>S. Mathlouthi Bouzid and C. Ben Othmane Zribi</i>	
Semantic Change Detection With Gaussian Word Embeddings . . . . .	3349
. . . . . <i>A. Yüksel, B. Uğurlu, and A. Koç</i>	
Medical Term and Status Generation From Chinese Clinical Dialogue With Multi-Granularity Transformer . . . . .	3362
. . . . . <i>M. Li, L. Xiang, X. Kang, Y. Zhao, Y. Zhou, and C. Zong</i>	
$F_0$ -Noise-Robust Glottal Source and Vocal Tract Analysis Based on ARX-LF Model . . . . .	3375
. . . . . <i>Y. Li, J. Tao, D. Erickson, B. Liu, and M. Akagi</i>	
Efficient Estimate of Sentence's Representation Based on the Difference Semantics Model . . . . .	3384
. . . . . <i>X. Liao, Y. Huang, Y. Wei, C. Zhang, F. Wang, and Y. Wang</i>	
TAU-Net: Temporal Activation U-Net Shared With Nonnegative Matrix Factorization for Speech Enhancement in Unseen Noise Environments . . . . .	3400
. . . . . <i>K. M. Jeon, G. W. Lee, N. K. Kim, and H. K. Kim</i>	
Robustness of Speech Spoofing Detectors Against Adversarial Post-Processing of Voice Conversion . . . . .	3415
. . . . . <i>Y.-Y. Ding, H.-J. Lin, L.-J. Liu, Z.-H. Ling, and Y. Hu</i>	
Language Agnostic Speaker Embedding for Cross-Lingual Personalized Speech Generation . . . . .	3427
. . . . . <i>Y. Zhou, X. Tian, and H. Li</i>	
Speech Enhancement Using Multi-Stage Self-Attentive Temporal Convolutional Networks . . . . .	3440
. . . . . <i>J. Lin, A. J. de L. van Wijngaarden, K.-C. Wang, and M. C. Smith</i>	

---

HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units . . . . .	3451
..... <i>W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed</i>	
Time-Frequency-Bin-Wise Linear Combination of Beamformers for Distortionless Signal Enhancement . . . . .	3461
..... <i>K. Yamaoka, N. Ono, and S. Makino</i>	
Convolutional Prediction for Monaural Speech Dereverberation and Noisy-Reverberant Speaker Separation . . . . .	3476
..... <i>Z.-Q. Wang, G. Wichern, and J. L. Roux</i>	
Learning Deep Direct-Path Relative Transfer Function for Binaural Sound Source Localization . . . . .	3491
..... <i>B. Yang, H. Liu, and X. Li</i>	
Pre-Training With Whole Word Masking for Chinese BERT . . . . .	3504
..... <i>Y. Cui, W. Che, T. Liu, B. Qin, and Z. Yang</i>	
Counterfactually Fair Automatic Speech Recognition . . . . .	3515
..... <i>L. Sari, M. Hasegawa-Johnson, and C. D. Yoo</i>	
Multi-Channel Multi-Frame ADL-MVDR for Target Speech Separation . . . . .	3526
..... <i>Z. Zhang, Y. Xu, M. Yu, S.-X. Zhang, L. Chen, D. S. Williamson, and D. Yu</i>	
Reduction of Subjective Listening Effort for TV Broadcast Signals With Recurrent Neural Networks . . . . .	3541
..... <i>N. L. Westhausen, R. Huber, H. Baumgartner, R. Sinha, J. RENNIES, and B. T. Meyer</i>	
Subword-Based Compact Reconstruction for Open-Vocabulary Neural Word Embeddings . . . . .	3551
..... <i>S. Sasaki, J. Suzuki, and K. Inui</i>	
Asynchronous Decentralized Distributed Training of Acoustic Models . . . . .	3565
..... <i>X. Cui, W. Zhang, A. Kayi, M. Liu, U. Finkler, B. Kingsbury, G. Saon, and D. Kung</i>	
Spatial Active Noise Control in Rooms Using Higher Order Sources . . . . .	3577
..... <i>J. Zhang, W. Zhang, J. A. Zhang, T. D. Abhayapala, and L. Zhang</i>	
Multimodal Emotion Recognition With Temporal and Semantic Consistency . . . . .	3592
..... <i>B. Chen, Q. Cao, M. Hou, Z. Zhang, G. Lu, and D. Zhang</i>	
Regularized Phrase-Based Topic Model for Automatic Question Classification With Domain-Agnostic Class Labels . . . . .	3604
..... <i>S. Supraja, A. W. H. Khong, and S. Tatinati</i>	
Sound Field Reproduction With a Cylindrical Loudspeaker Array Using First Order Wall Reflections . . . . .	3617
..... <i>N. Maeda, F. M. Fazi, and F.-M. Hoffmann</i>	
Coupling a Generative Model With a Discriminative Learning Framework for Speaker Verification . . . . .	3631
..... <i>X. Lu, P. Shen, Y. Tsao, and H. Kawai</i>	
Effects of Additive Noise in Binaural Rendering of Spherical Microphone Array Signals . . . . .	3642
..... <i>H. Helmholtz, D. Lou Alon, S. V. A. Garí, and J. Ahrens</i>	
Speech Reconstruction With Reminiscent Sound Via Visual Voice Memory . . . . .	3654
..... <i>J. Hong, M. Kim, S. J. Park, and Y. M. Ro</i>	
Robustness of Acoustic Rake Filters in Minimum Variance Beamforming . . . . .	3668
..... <i>R. Weisman, T. Shlomo, V. Tourbabin, P. Calamia, and B. Rafaely</i>	
Mixed Precision Low-Bit Quantization of Neural Network Language Models for Speech Recognition . . . . .	3679
..... <i>J. Xu, J. Yu, S. Hu, X. Liu, and H. Meng</i>	
Learning Fine-Grained Fact-Article Correspondence in Legal Cases . . . . .	3694
..... <i>J. Ge, Y. Huang, X. Shen, C. Li, and W. Hu</i>	
High-Resolution Piano Transcription With Pedals by Regressing Onset and Offset Times . . . . .	3707
..... <i>Q. Kong, B. Li, X. Song, Y. Wan, and Y. Wang</i>	

---