# 26th Annual Conference of the International Speech Communication Association (INTERSPEECH 2025)

Rotterdam, The Netherlands
17-21 August 2025

Volume 1 of 8

# TABLE OF CONTENTS

## VOLUME 1

## MULTILINGUALITY, CROSS-LINGUISTIC STUDIES, L2 SPEECH

## SPEECH EMOTION RECOGNITION 1

## MULTIMODAL RESOURCES

# INTERPRETABILITY IN AUDIO AND SPEECH TECHNOLOGY

## SUMMARIZATION

## SHOW AND TELL 1: ASR / TOOLS

## MODELS OF SPEECH PRODUCTION

## SPEECH AND GRAMMAR/ARTICULATORY ANALYSES

## SPEAKING STYLES, REGISTER AND CONVERSATIONAL SPEECH

# EMOTIONAL DISTRESS IN SPEECH

# PROSODY IN SPEECH SYNTHESIS

## DEPRESSION DETECTION AND ASSESSMENT 1

## SPEECH ANALYSIS, DETECTION AND CLASSIFICATION 1

## SPEECH-BASED COGNITIVE ASSESSMENT 1

## LARGE LANGUAGE MODELS IN SPEECH RECOGNITION

# SPEECH CODING AND ECHO CANCELLATION

# VOLUME 2

# DECODING ALGORITHMS

## QUEER AND TRANS SPEECH SCIENCE AND TECHNOLOGY

## TONE

## CROSS-LINGUAL AND MULTILINGUAL PROCESSING

## ECHO CANCELLATION, FEEDBACK CONTROL, AND NEAR-END ENHANCEMENT

## PATHOLOGICAL SPEECH ANALYSIS 1

## HEARING DISORDERS

## INTERSPEECH 2025 URGENT CHALLENGE

## SPOKEN MACHINE TRANSLATION 2

## SPATIAL AUDIO AND ACOUSTICS 1

## ARTICULATORY AND VOCAL TRACT MODELLING

## ACOUSTIC ASSESSMENT OF RESPIRATORY HEALTH

## ADVANCES IN MODELLING AND IMAGING

## CONVERSATION, COMMUNICATION AND INTERACTION 1

## ROBUST SPEAKER VERIFICATION

## MULTILINGUAL ASR

## MULTI-CHANNEL SPEECH ENHANCEMENT

## SELF-SUPERVISED LEARNING

## SINGING VOICE AND AUDIO SYNTHESIS

## ACOUSTIC AND ARTICULATORY CUES IN SPEECH PERCEPTION

## AUDIO EVENT DETECTION AND CLASSIFICATION

# INCLUSIVITY

# VOICE CONVERSION 1

# VOLUME 3

## SPEECH-BASED COGNITIVE ASSESSMENT 2

## SOURCE SEPARATION 1

## LANGUAGE AND ACCENT IDENTIFICATION AND SPEAKER PRIVACY

## SOURCE TRACING: THE ORIGINS OF SYNTHETIC OR MANIPULATED SPEECH

## SPEAKER DIARIZATION 1

## MULTILINGUAL SPEECH SYNTHESIS AND SPECIAL APPLICATIONS 1

## CHARACTERIZATION AND MULTIMODAL APPROACHES FOR SPEAKER RECOGNITION

## ACOUSTIC ANALYSIS AND BIOACOUSTICS

## SPOKEN DIALOGUE SYSTEMS 1

## SPEECH ASSESSMENT

## AUDIO-VISUAL ASR AND MULTIMODAL SYSTEM

## SPEECH AND VOICE DISORDERS 1

## MULTIMODAL INFORMATION BASED SPEECH PROCESSING (MISP) 2025 CHALLENGE

## SPEAKER EXTRACTION 1

## LOW RESOURCE SPEECH RECOGNITION

## COMPUTATIONAL RESOURCE CONSTRAINED ASR

## SPEECH AND LANGUAGE TECHNOLOGY FOR HEALTH APPLICATIONS

## RESPONSIBLE SPEECH FOUNDATION MODELS + SUPERB CHALLENGE

## DYSARTHRIC SPEECH ASSESSMENT 1

## VOLUME 4

## SHOW AND TELL 2: SPEECH SYNTHESIS

## DATABASES AND PROGRESS IN METHODOLOGY

## NOVEL ARCHITECTURES FOR ASR

## DEEPFAKE DETECTION

## TOOLS FOR SPEECH ANALYSIS

## TEXT PROCESSING AND EVALUATION FOR SPEECH SYNTHESIS 1

## SEGMENTAL AND TONAL UNITS

## SPEECH QUALITY ASSESSMENT

## SPEECH ENHANCEMENT

# LANGUAGE LEARNING AND ASSESSMENT

# SPEECH SYNTHESIS PARADIGMS AND METHODS 1

## SPATIAL AUDIO AND ACOUSTICS 2

# TEXT PROCESSING AND EVALUATION FOR SPEECH SYNTHESIS 2

# GENERAL TOPICS IN ASR

# ACOUSTIC EVENT DETECTION AND CLASSIFICATION

## KEYWORD SPOTTING AND RETRIEVAL

## **MULTIMODAL SYSTEMS**

# DYSARTHRIC SPEECH ASSESSMENT 2

# DIALECT IDENTIFICATION IN DIFFERENT LANGUAGES

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

# VOLUME 5

## BRAIN AND COGNITION

## REGIONAL, SOCIAL AND DIACHRONIC VARIATION

## SPEAKER EXTRACTION 2

## MULTIMODAL EMOTION RECOGNITION

## CONVERSATION, COMMUNICATION AND INTERACTION 2

## MULTIMODAL SPEECH AND LANGUAGE PROCESSING IN HEALTHCARE SETTINGS

## MUSIC AND AUDIO ANALYSIS

## AUDIO ANALYSIS, GENERATION AND ASSESSMENT

## OTHER TOPICS IN SPEECH RECOGNITION

# PRIVACY AND ANONYMIZATION

# LANGUAGE MODELING FOR CONVERSATIONAL SYSTEMS

# SPEECH ACCESSIBILITY PROJECT CHALLENGE

## NEURAL NETWORK TRAINING METHODS 1

## DIVERSITY: AGE, SEX, GENDER, ETHNICITY, AND MORE

## ANOMALOUS SOUND DETECTION

## FAR-FIELD AND ROBUST SPEECH RECOGNITION

## SPEECH SYNTHESIS PARADIGMS AND METHODS 2

## ARTICULATORY ANALYSES

## SPEECH AND AUDIO ANALYSIS AND REPRESENTATION

## SHOW AND TELL 3: SIGNAL PROCESSING / MULTIMODAL PROCESSING

## SPEECH AND VOICE DISORDERS 2

## NEURAL NETWORK TRAINING METHODS 2

# VOLUME 6

## TRAINING AND SCORING METHODS FOR SPEAKER RECOGNITION

## PATHOLOGICAL SPEECH ANALYSIS 2

## MULTIMODAL AND VISUAL SPEECH SYNTHESIS

## LEXICON AND GRAMMAR

## NOISE REDUCTION AND DEREVERBERATION

## NEURAL NETWORK TRAINING METHODS AND ARCHITECTURES

## CHALLENGES IN SPEECH DATA COLLECTION, CURATION AND ANNOTATION - PART 1

## EVALUATION AND FORENSIC APPLICATIONS OF SPEAKER RECOGNITION

## LANGUAGE RESOURCES

# BANDWIDTH EXPANSION AND DIFFUSION-BASED SPEECH ENHANCEMENT

# SPOKEN LANGUAGE UNDERSTANDING

# MULTILINGUAL SPEECH SYNTHESIS AND SPECIAL APPLICATIONS 2

# PROSODY AND VOICE QUALITY

## GENERATIVE MODELS FOR AUDIO

## CHALLENGES IN SPEECH DATA COLLECTION, CURATION AND ANNOTATION - PART 2

## SPEECH EMOTION RECOGNITION 3

# VOLUME 7

## EMOTION AND EXPRESSIVITY IN SPEECH SYNTHESIS AND VOICE CONVERSION

## STREAMING ASR

## L1 AND L2 ACQUISITION, PERCEPTION AND PROCESSING

## SPEECH EMOTION RECOGNITION 2

## SPEAKER TRAITS RECOGNITION

## SPOOFING AND ADVERSARIAL ATTACKS

## VOICE CONVERSION 2

# PATHOLOGICAL SPEECH ANALYSIS 3

# SPEECH EMOTION RECOGNITION IN NATURALISTIC CONDITIONS CHALLENGE

## PROSODY, PHONEME AND STRESS MODELING IN ASR

## SEGMENTS

## DATASETS AND TOOLS FOR SPEECH SYNTHESIS

## SPOKEN DIALOGUE SYSTEMS 2

## SPEECH ENHANCEMENT AND REPRESENTATION LEARNING

## NEURAL CODECS AND VOCODERS

## ADAPTATION AND TARGET-SPEAKER ASR

## SHOW AND TELL 4: EDUCATION / ASSISTIVE TECHNOLOGY

## SOURCE SEPARATION 2

## SPEECH CODING

## MULTIMODALITY

## SPEECH ASSESSMENT AND LANGUAGE LEARNING

## VOLUME 8

## WATERMARKING AND ANONYMIZATION

## SINGLE-CHANNEL SPEECH ENHANCEMENT

## CONTEXTUAL BIASING AND ADAPTATION

## SPEAKER DIARIZATION 2

## DEPRESSION DETECTION AND ASSESSMENT 2

# PATHOLOGICAL SPEECH ANALYSIS 4

# SPEECH DEEPFAKES

## PROSODY

## SPEECH ANALYSIS AND QUALITY ASSESSMENT

## EMOTIONS AND FOUNDATIONAL MODELS

## PREDICTION AND EVALUATION OF SPEECH QUALITY AND INTELLIGIBILITY

## MULTI-TALKER ASR

## SPEECH SYNTHESIS PARADIGMS AND METHODS 3

# BIOSIGNAL-ENABLED SPOKEN COMMUNICATION

## SPEECH DEEPFAKES, ANTISPOOFING AND BACKDOOR ATTACKS

## PATHOLOGICAL SPEECH ANALYSIS 5

## ASR ASSESSMENT AND FOUNDATIONAL MODELS

## SPEAKER RECOGNITION

## SPEECH ANALYSIS, DETECTION AND CLASSIFICATION 2

**Author Index**