

IEEE SLT 2024 Program Details

Day 1, Dec 2, Monday

08:30-09:00 Opening Session (Venue: Lecture Hall)

09:00-10:00 Keynote Speech 1 (Venue: Lecture Hall)

Title: Towards Robust Audio Deepfake Detection and Attribution

Speaker: *Prof Jianhua Tao, Tsinghua University*

Chair: *Dr Minghui Dong*

10:00-10:30 Coffee Break

10:30-12:30 Poster Session 1: Speech Recognition (Venue: Poster Area)

Chair: Prof Yanmin Qian

Poster ID (Paper ID)	Title and Authors
P1-01-ASR (#28)	PromptKWS: A Novel Prompt-Guided Open-Vocabulary Keyword Spotting Framework <i>Gaopeng Xu (NIO)</i> <i>Chengfei Li (Qilu Normal University)</i> <i>Xianliang Wang (NIO)</i> <i>Li Zhu (NIO)</i> <i>Juan Wei (NIO)</i> <i>Wenpeng Li (NIO)</i> <i>Jianwei Niu (NIO)</i> <i>Jie Gao (NIO)</i>
P1-02-ASR (#43)	Personalizing Large Sequence-to-Sequence Speech Foundation Models with Speaker Representations <i>Dominik Wagner (Technische Hochschule Nürnberg Georg Simon Ohm)</i> <i>Ilja Baumann (Technische Hochschule Nürnberg Georg Simon Ohm)</i> <i>Thomas Ranzenberger (Technische Hochschule Nürnberg Georg Simon Ohm)</i> <i>Korbinian Riedhammer (Technische Hochschule Nürnberg Georg Simon Ohm)</i> <i>Tobias Bocklet (TH Nürnberg)</i>

P1-03-ASR (#97)	<p>Label-Looping: Highly Efficient Decoding for Transducers</p> <p><i>Vladimir Bataev (NVIDIA, University of London)</i> <i>Hainan Xu (NVIDIA)</i> <i>Daniel Galvez (NVIDIA)</i> <i>Vitaly Lavrukhin (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i></p>
P1-04-ASR (#102)	<p>Advancing Multi-Talker ASR Performance with Large Language Models</p> <p><i>Mohan Shi (University of California, Los Angeles)</i> <i>Zengrui Jin (The Chinese University of Hong Kong)</i> <i>Yaoxun Xu (Tsinghua University)</i> <i>Yong Xu (Tencent)</i> <i>Shi-Xiong Zhang (Capital One)</i> <i>Kun Wei (School of Computer Science, Northwestern Polytechnical University)</i> <i>Yiwen Shao (Johns Hopkins University)</i> <i>Chunlei Zhang (Bytedance)</i> <i>Dong Yu (Tencent AI Lab)</i></p>
P1-05-ASR (#114)	<p>Token-Weighted RNN-T for Learning from Flawed Data</p> <p><i>Gil Keren (Meta)</i> <i>Wei Zhou (Meta)</i> <i>Ozlem Kalinli (Meta)</i></p>
P1-06-ASR (#137)	<p>Enhancing Code-Switching Speech Recognition with LID-Based Collaborative Mixture of Experts Model</p> <p><i>Hukai Huang (Xiamen University)</i> <i>Jiayan Lin (Xiamen University)</i> <i>Kaidi Wang (Xiamen University)</i> <i>Yishuang Li (Xiamen University)</i> <i>Wenhao Guan (Xiamen University)</i> <i>Lin Li (Xiamen University)</i> <i>Qingyang Hong (Xiamen University)</i></p>
P1-07-ASR (#149)	<p>Language Bias in Self-Supervised Learning for Automatic Speech Recognition</p> <p><i>Ed Storey (Trinity College Dublin)</i> <i>Naomi Harte (Trinity College Dublin)</i> <i>Peter Bell (University of Edinburgh)</i></p>
P1-08-ASR (#150)	<p>Robust Audiovisual Speech Recognition Models with Mixture-of-Experts</p>

	<p><i>Yihan Wu (Renmin University of China)</i> <i>Yifan Peng (Carnegie Mellon University)</i> <i>Yichen Lu (Carnegie Mellon University)</i> <i>Xuankai Chang (Carnegie Mellon University)</i> <i>Ruihua Song (Renmin University of China)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i></p>
P1-09-ASR (#162)	<p>Hybrid Attention-Based Encoder-Decoder Model for Efficient Language Model Adaptation</p> <p><i>Shaoshi Ling (Microsoft)</i> <i>Guoli Ye (Microsoft)</i> <i>Rui Zhao (Microsoft)</i> <i>Yifan Gong (Microsoft)</i></p>
P1-10-ASR (#165)	<p>SpatialEmb: Extract and Encode Spatial Information for 1-Stage Multi-Channel Multi-Speaker ASR on Arbitrary Microphone Arrays</p> <p><i>Yiwen Shao (Johns Hopkins University)</i> <i>Yong Xu (Tencent)</i> <i>Sanjeev Khudanpur (Johns Hopkins University)</i> <i>Dong Yu (Tencent AI Lab)</i></p>
P1-11-ASR (#171)	<p>Effective Text Adaptation for LLM-Based ASR through Soft Prompt Fine-Tuning</p> <p><i>Yingyi Ma (Meta)</i> <i>Zhe Liu (Meta)</i> <i>Ozlem Kalinli (Meta)</i></p>
P1-12-ASR (#173)	<p>Temporal Order Preserved Optimal Transport-Based Cross-Modal Knowledge Transfer Learning for ASR</p> <p><i>Xugang Lu (NICT)</i> <i>Peng Shen (NICT)</i> <i>Yu Tsao (Academia Sinica)</i> <i>Hisashi Kawai (NICT)</i></p>
P1-13-ASR (#187)	<p>Contextualized Automatic Speech Recognition with Dynamic Vocabulary</p> <p><i>Yui Sudo (Honda Research Institute Japan)</i> <i>Yosuke Fukumoto (Honda Research Institute Japan)</i> <i>Muhammad Shakeel (Honda Research Institute Japan)</i> <i>Yifan Peng (Carnegie Mellon University)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i></p>
P1-14-ASR	<p>Do Prompts Really Prompt? Exploring the Prompt Understanding</p>

(#209)	<p>Capability of Whisper</p> <p><i>Chih-Kai Yang (National Taiwan University)</i> <i>Kuan-Po Huang (National Taiwan University)</i> <i>Hung-yi Lee (National Taiwan University)</i></p>
P1-15-ASR (#221)	<p>An Effective Context-Balanced Adaptation Approach for Long-Tailed Speech Recognition</p> <p><i>Yi-Cheng Wang (National Taiwan Normal University)</i> <i>Li-Ting Pai (National Taiwan Normal University)</i> <i>Bi-Cheng Yan (National Taiwan Normal University)</i> <i>Hsin-Wei Wang (NTNU)</i> <i>Chi-Han Lin (E.SUN Financial Holding Co., Ltd.)</i> <i>Berlin Chen (National Taiwan Normal University)</i></p>
P1-16-ASR (#300)	<p>Training Large ASR Encoders with Differential Privacy</p> <p><i>Geeticka Chauhan (Google DeepMind)</i> <i>Steve Chien (Google)</i> <i>Om Thakkar (Google)</i> <i>Abhradeep Thakurta (Google)</i> <i>Arun Narayanan (Google Inc.)</i></p>
P1-17-ASR (#309)	<p>Transducer Consistency Regularization for Speech-to-Text Applications</p> <p><i>Cindy S Tseng (Samsung Research America)</i> <i>Yun Tang (Samsung Research America)</i> <i>Vijendra Raj Apsingekar (Samsung Research America)</i></p>
P1-18-ASR (#324)	<p>Leave No Knowledge Behind During Knowledge Distillation: Practical and Effective Knowledge Distillation for Code-Switching ASR Using Realistic Data</p> <p><i>Liang-Hsuan Tseng (National Taiwan University)</i> <i>Zih-Ching Chen (National Taiwan University)</i> <i>Weishun Chang (National Taiwan University)</i> <i>Cheng-Kuang Lee (NVIDIA Corporation)</i> <i>Tsung-Ren Huang (National Taiwan University)</i> <i>Hung-yi Lee (National Taiwan University)</i></p>
P1-19-ASR (#389)	<p>CTC-Assisted LLM-Based Contextual ASR</p> <p><i>Guanrou Yang (Shanghai Jiao Tong University)</i> <i>Ziyang Ma (Shanghai Jiao Tong University)</i> <i>Zhifu Gao (Alibaba)</i> <i>Shiliang Zhang (Alibaba Group)</i> <i>Xie Chen (Shanghai Jiao Tong University)</i></p>

P1-20-ASR (#53)	<p>Automatic Time Alignment Generation for End-to-End ASR Using Acoustic Probability Modeling</p> <p><i>Dongcheng Jiang (University of Cambridge)</i> <i>Chao Zhang (Tsinghua University)</i> <i>Phil Woodland (Machine Intelligence Laboratory, Cambridge University Department of Engineering)</i></p>
P1-21-ASR (#73)	<p>Continual Learning with Embedding Layer Surgery and Task-Wise Beam</p> <p><i>Chin Yuen Kwok (Nanyang Technological University)</i> <i>Jia Qi Yip (Alibaba Group / Nanyang Technological University)</i> <i>Eng Siong Chng (Nanyang Technological University)</i></p>
P1-22-ASR (#93)	<p>BESTOW: Efficient and Streamable Speech Language Model with the Best of GPT and T5</p> <p><i>He Huang (NVIDIA)</i> <i>Zhehuai Chen (NVIDIA)</i> <i>Krishna C Puvvada (NVIDIA)</i> <i>Piotr Żelasko (NVIDIA)</i> <i>Jagadeesh Balam (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i> <i>Nithin Rao Koluguri (NVIDIA)</i> <i>Oleksii Hrinchuk (NVIDIA)</i></p>
P1-23-ASR (#116)	<p>Combining TF-GridNet and Mixture Encoder for Continuous Speech Separation for Meeting Transcription</p> <p><i>Peter Vieting (RWTH Aachen University)</i> <i>Simon Berger (RWTH Aachen University)</i> <i>Thilo von Neumann (Paderborn University)</i> <i>Christoph Boeddeker (Paderborn University)</i> <i>Ralf Schlüter (RWTH Aachen University)</i> <i>Reinhold Haeb-Umbach (Paderborn University)</i></p>
P1-24-ASR (#122)	<p>Mamba-Based Decoder-Only Approach with Bidirectional Speech Modeling for Speech Recognition</p> <p><i>Yoshiki Masuyama (Tokyo Metropolitan University)</i> <i>Koichi Miyazaki (CyberAgent)</i> <i>Masato Murata (CyberAgent)</i></p>
P1-25-ASR (#201)	<p>An Analysis of Linear Complexity Attention Substitutes with BEST-RQ</p> <p><i>Ryan Whetten (LIA - Avignon University)</i></p>

	<i>Titouan Parcollet (Samsung AI Cambridge / University of Cambridge)</i> <i>Adel Moumen (Avignon University)</i> <i>Marco Dinarelli (CNRS)</i> <i>Yannick Estève (LIA - Avignon University)</i>
P1-26-ASR (#214)	Speech-Mamba: Long-Context Speech Recognition with Selective State Spaces Models <i>Xiaoxue Gao (ASTAR)</i> <i>Nancy Chen (Institute for Infocomm Research)</i>
P1-27-ASR (#357)	Lite ASR Transformer: A Lightweight Transformer Architecture for Automatic Speech Recognition <i>Metilda Sagaya Mary NJ (Indian Institute of Technology Madras)</i> <i>S Umesh (IIT Chennai)</i>

10:30-12:30 Challenge Session 1: Stutter Speech ASR/SED and Dysarthria WWS (Venue: Lecture Hall)

12:30-14:00 Lunch

14:00-15:00 Invited Talk 1 (Venue: Lecture Hall)

Title: Challenges and Progress in Automatic Speech-to-Speech Translation: Bridging the Gap to Real-Time Interpretation

Speaker: *Prof Satoshi Nakamura, The Chinese University of Hong Kong, Shenzhen*

Chair: *Prof Xie Chen*

15:00-15:30 Coffee Break

15:30-17:30 Poster Session 2: Speech Recognition and Enhancement (Venue: Poster Area)

Chair: Prof Jun Du

Poster ID (Paper ID)	Title and Authors
P2-01-ASR (#22)	Serialized Speech Information Guidance with Overlapped Encoding Separation for Multi-Speaker Automatic Speech Recognition <i>Hao Shi (Kyoto University)</i> <i>Yuan Gao (Kyoto University)</i> <i>Zhaoheng Ni (Meta AI)</i> <i>Tatsuya Kawahara (Kyoto University)</i>
P2-02-ASR	Efficient Extraction of Noise-Robust Discrete Units from Self-

(#161)	<p>Supervised Speech Models</p> <p><i>Jakob Poncelet (KU Leuven)</i> <i>Yujun Wang (Xiaomi)</i> <i>Hugo Van Hamme (KU Leuven)</i></p>
P2-03-ASR (#170)	<p>Controlling Whisper: Universal Acoustic Adversarial Attacks to Control Multi-Task Automatic Speech Recognition Models</p> <p><i>Vyas Raina (University of Cambridge)</i> <i>Mark Gales (University of Cambridge)</i></p>
P2-04-ASR (#283)	<p>Improving Rare-Word Recognition of Whisper in Zero-Shot Settings</p> <p><i>Yash Jogi (Sprinklr)</i> <i>Vaibhav Aggarwal (Sprinklr)</i> <i>Shabari S Nair (Sprinklr)</i> <i>Yash Verma (Sprinklr)</i> <i>Aayush Kubba (Sprinklr)</i></p>
P2-05-ASR (#343)	<p>Augmenting Automatic Speech Recognition Models with Disfluency Detection</p> <p><i>Robin Amann (Karlsruher Institut für Technologie)</i> <i>Zhaolin Li (Karlsruhe Institute of Technology)</i> <i>Barbara Bruno (Karlsruhe Institute of Technology)</i> <i>Jan Niehues (Karlsruhe Institute of Technology)</i></p>
P2-06-ASR (#246)	<p>Enhancing Unified Streaming and Non-Streaming ASR through Curriculum Learning with Easy-to-Hard Tasks</p> <p><i>Yuting Yang (NetEase Yidun AI Lab)</i> <i>Yuke Li (NetEase Yidun AI Lab)</i> <i>Lifeng Zhou (NetEase Yidun AI Lab)</i> <i>Binbin Du (NetEase Yidun AI Lab)</i> <i>Haoqi Zhu (NetEase Yidun AI Lab)</i></p>
P2-07-ASR (#78)	<p>DQ-Whisper: Joint Distillation and Quantization for Efficient Multilingual Speech Recognition</p> <p><i>Hang Shao (Shanghai Jiao Tong University)</i> <i>Bei Liu (Shanghai Jiao Tong University)</i> <i>Wei Wang (Shanghai Jiao Tong University)</i> <i>Xun Gong (Shanghai Jiao Tong University)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i></p>
P2-08-ASR (#129)	<p>Fusion of Discrete Representations and Self-Augmented Representations for Multilingual Automatic Speech Recognition</p>

	<p><i>Shih-Heng Wang (National Taiwan University)</i> <i>Jiatong Shi (Carnegie Mellon University)</i> <i>Chien-Yu Huang (National Taiwan University)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P2-09-ASR (#157)	<p>Longer is (Not Necessarily)</p> <p><i>Nithin Rao Koluguri (NVIDIA)</i> <i>Travis M Bartley (NVIDIA CUNY)</i> <i>Hainan Xu (NVIDIA)</i> <i>Oleksii Hrinchuk (NVIDIA)</i> <i>Jagadeesh Balam (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i> <i>Georg Kucsko (Suno Inc.)</i></p>
P2-10-ASR (#178)	<p>Semi-Supervised Learning for Code-Switching ASR with Large Language Model Filter</p> <p><i>Yu Xi (NVIDIA)</i> <i>Wen Ding (NVIDIA)</i> <i>Kai Yu (Shanghai Jiao Tong University)</i> <i>Junjie Lai (NVIDIA)</i></p>
P2-11-ASR (#301)	<p>Parameter Averaging is All You Need to Prevent Forgetting</p> <p><i>Peter W Plantinga (JPMorgan Chase & Co.)</i> <i>Jaekwon Yoo (JPMorgan Chase & Co.)</i> <i>Abenezer G Girma (JP Morgan Chase)</i> <i>Chandra Dhir (JPMorgan Chase)</i></p>
P2-12-ASR (#304)	<p>Advancing CTC Models for Better Speech Alignment: A Topological Approach</p> <p><i>Zeyu Zhao (University of Edinburgh)</i> <i>Peter Bell (University of Edinburgh)</i></p>
P2-13-SES (#37)	<p>DualSep: A Lightweight Dual-Encoder Convolutional Recurrent Network for Real-Time In-Car Speech Separation</p> <p><i>Ziqian Wang (Northwestern Polytechnical University)</i> <i>Jiayao Sun (Northwestern Polytechnical University)</i> <i>Zihan Zhang (Northwestern Polytechnical University)</i> <i>Xingchen Li (Northwestern Polytechnical University)</i> <i>Jie Liu (Huawei Cloud)</i> <i>Lei Xie (NWPU)</i></p>

<p>P2-14-SES (#39)</p>	<p>DDTSE: Discriminative Diffusion Model for Target Speech Extraction</p> <p><i>Leying Zhang (Shanghai Jiao Tong University)</i> <i>Yao Qian (Microsoft)</i> <i>Linfeng Yu (Shanghai Jiao Tong University)</i> <i>Heming Wang (The Ohio State University)</i> <i>Hemin Yang (Microsoft)</i> <i>Shujie Liu (Microsoft Research Asia)</i> <i>Long Zhou (Microsoft Research Asia)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i></p>
<p>P2-15-SES (#89)</p>	<p>An Investigation of Incorporating Mamba for Speech Enhancement</p> <p><i>Rong Chao (National Taiwan University)</i> <i>Wen-Huang Cheng (National Taiwan University)</i> <i>Moreno La Quatra (Kore University of Enna)</i> <i>Sabato M Siniscalchi (University of Palermo)</i> <i>Chao-Han Huck Yang (NVIDIA Research)</i> <i>Szu-Wei Fu (NVIDIA)</i> <i>Yu Tsao (Academia Sinica)</i></p>
<p>P2-16-SES (#117)</p>	<p>Effective Noise-Aware Data Simulation for Domain-Adaptive Speech Enhancement Leveraging Dynamic Stochastic Perturbation</p> <p><i>Chien-Chun Wang (National Taiwan Normal University)</i> <i>Li-Wei Chen (United Link Co., Ltd.)</i> <i>Hung-Shin Lee (United Link Co., Ltd.)</i> <i>Berlin Chen (National Taiwan Normal University)</i> <i>Hsin-Min Wang (Academia Sinica)</i></p>
<p>P2-17-SES (#120)</p>	<p>SMRU: Split-and-Merge Recurrent-Based UNet for Acoustic Echo Cancellation and Noise Suppression</p> <p><i>Zhihang Sun (Tencent AI Lab)</i> <i>Andong Li (Tencent AI Lab)</i> <i>Rilin Chen (Tencent)</i> <i>Hao Zhang (Tencent AI Lab)</i> <i>Meng Yu (Tencent)</i> <i>Yi Zhou (CQUPT)</i> <i>Dong Yu (Tencent AI Lab)</i></p>
<p>P2-18-SES (#139)</p>	<p>On the Effectiveness of Enrollment Speech Augmentation for Target Speaker Extraction</p> <p><i>Junjie Li (The Hong Kong Polytechnic University)</i> <i>Ke Zhang (Northeastern University)</i> <i>Shuai Wang (Shenzhen Research Institute of Big Data, Chinese University)</i></p>

	<p><i>of Hong Kong (Shenzhen)</i> <i>Haizhou Li (The Chinese University of Hong Kong, Shenzhen)</i> <i>M W Mak (HK PolyU)</i> <i>Kong Aik Lee (The Hong Kong Polytechnic University)</i></p>
P2-19-SES (#142)	<p>Diffusion-Based Generative Modeling with Discriminative Guidance for Streamable Speech Enhancement</p> <p><i>Chenda Li (Shanghai Jiao Tong University)</i> <i>Samuele Cornell (Carnegie Mellon University)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i></p>
P2-20-SES (#216)	<p>NeuroSpex: Neuro-Guided Speaker Extraction with Cross-Modal Fusion</p> <p><i>Dashanka D N De Silva (University of Bremen)</i> <i>Siqi Cai (National University of Singapore)</i> <i>Saurav Pahuja (University of Bremen)</i> <i>Tanja Schultz (University of Bremen)</i> <i>Haizhou Li (The Chinese University of Hong Kong, Shenzhen)</i></p>
P2-21-SES (#325)	<p>Enhancing Speaker Extraction through Rectifying Target Confusion</p> <p><i>Jiahe Wang (Shanghai Jiao Tong University)</i> <i>Shuai Wang (Shenzhen Research Institute of Big Data, Chinese University of Hong Kong (Shenzhen))</i> <i>Junjie Li (The Hong Kong Polytechnic University)</i> <i>Ke Zhang (Northeastern University)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i> <i>Haizhou Li (The Chinese University of Hong Kong, Shenzhen)</i></p>
P2-22-SES (#368)	<p>Diff-PLC: A Diffusion-Based Approach for Effective Packet Loss Concealment</p> <p><i>Da-Hee Yang (Hanyang University)</i> <i>Joon-Hyuk Chang (Hanyang University)</i></p>
P2-23-SES (#369)	<p>Improving Curriculum Learning for Target Speaker Extraction with Synthetic Speakers</p> <p><i>Yun Liu (National Institute of Informatics)</i> <i>Xuechen Liu (National Institute of Informatics)</i> <i>Junichi Yamagishi (National Institute of Informatics)</i></p>
P2-24-SS02 (#415)	<p>Language Model Based Generative Error Correction: A Challenge and Baselines for Speech Recognition, Speaker Tagging, and Emotion Recognition</p>

	<p> <i>Chao-Han Huck Yang (NVIDIA Research)</i> <i>Tae Jin Park (NVIDIA)</i> <i>Yuan Gong (Massachusetts Institute of Technology)</i> <i>Yuanchao Li (University of Edinburgh)</i> <i>Yen-Ting Lin (National Taiwan University)</i> <i>Zhehuai Chen (NVIDIA)</i> <i>Yuchen Hu (Nanyang Technological University)</i> <i>Chen Chen (Nanyang Technological University)</i> <i>Kunal Dhawan (NVIDIA)</i> <i>Piotr Żelasko (NVIDIA)</i> <i>Chao Zhang (Tsinghua University)</i> <i>Yun-Nung Chen (National Taiwan University)</i> <i>Yu Tsao (Academia Sinica)</i> <i>Jagadeesh Balam (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i> <i>Sabato M Siniscalchi (University of Palermo)</i> <i>Eng Siong Chng (Nanyang Technological University)</i> <i>Peter Bell (University of Edinburgh)</i> <i>Catherine Lai (University of Edinburgh)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> <i>Andreas Stolcke (Uniphore)</i> </p>
P2-25-SS06 (#400)	<p> FGCL: Fine-Grained Contrastive Learning for Mandarin Stuttering Event Detection </p> <p> <i>Han Jiang (Xi'an Jiaotong University)</i> <i>Wenyu Wang (Xi'an Jiaotong University)</i> <i>Yiquan Zhou (XJTU)</i> <i>Hongwu Ding (Happy Elements)</i> <i>Xu Jiacheng (Happy Elements)</i> <i>Jihua Zhu (Xi'an Jiaotong University)</i> </p>
P2-26-SS06 (#402)	<p> Data Augmentation Techniques for Improved Performance in the SLT 2024 Mandarin Stuttering Event Detection and ASR Challenge </p> <p> <i>Weiwei Wang (Chery HuiYin Motor Finance Service Co., Ltd.)</i> <i>Zhijin Feng (Chery HuiYin Motor Finance Service Co., Ltd.)</i> <i>Qingyuan Song (Chery HuiYin Motor Finance Service Co., Ltd.)</i> <i>Wenyang Wei (Chery HuiYin Motor Finance Service Co., Ltd.)</i> <i>Yansong Wang (Chery HuiYin Motor Finance Service Co., Ltd.)</i> </p>
P2-27-SS06 (#403)	<p> Findings of the 2024 Mandarin Stuttering Event Detection and Automatic Speech Recognition Challenge </p> <p> <i>Hongfei Xue (NWPU)</i> <i>Rong Gong (StammerTalk)</i> <i>Mingchen Shao (NWPU)</i> <i>Xin Xu (AISHELL)</i> <i>Lezhi Wang (StammerTalk)</i> </p>

	<p><i>Lei Xie (NWPU)</i> <i>Hui Bu (AISHELL)</i> <i>Jiaming Zhou (Nankai University)</i> <i>Yong Qin (Nankai University)</i> <i>Jun Du (University of Science and Technology of China)</i> <i>Ming Li (Wuhan University)</i> <i>Binbin Zhang (WeNet Open Source Community)</i> <i>Bin Jia (StammerTalk)</i></p>
P2-28-SS06 (#410)	<p>Enhanced ASR for Stuttering Speech: Combining Adversarial and Signal-Based Data Augmentation</p> <p><i>Shangkun Huang (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Dejun Zhang (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Jing Deng (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Rong Zheng (Beijing Fosafer Information Technology Co., Ltd.)</i></p>

15:30-17:30 Challenge Session 2: Singing Voice Deepfake Detection (SVDD)
(Venue: Lecture Hall)

17:30-18:30 Recent Breakthrough (Venue: Poster Area)

Chair: Prof Eng Siong Chng

Poster ID (Paper ID)	Title and Authors
RB-1 (#1)	<p>Reverb: Open-Source ASR and Diarization from Rev</p> <p><i>Nishchal Bhandari, Danny Chen, Miguel Ángel del Río Fernández, Natalie Delworth, Jennifer Drexler Fox, Migüel Jetté, Quinten McNamara, Corey Miller, Ondřej Novotný, Ján Profant, Nan Qin, Martin Ratajczak, Jean-Philippe Robichaud</i></p>
RB-2 (#2)	<p>GAMA: A Large Audio-Language Model with Advanced Audio Understanding and Complex Reasoning Abilities</p> <p><i>Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Kiran Reddy Evuru, Utkarsh Tyagi, S Sakshi, Oriol Nieto, Ramani Duraiswami, Dinesh Manocha</i></p>
RB-3 (#3)	<p>MMAU: A Massive Multi-Task Audio Understanding and Reasoning Benchmark</p> <p><i>S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, Dinesh Manocha</i></p>
RB-4 (#4)	<p>Vevo: Controllable Zero-Shot Voice Imitation with Self-Supervised Disentanglement</p>

	<i>Xueyao Zhang, Xiaohui Zhang, Kainan Peng, Zhenyu Tang, Vimal Manohar, Yingru Liu, Jeff Hwang, Dangna Li, Yuhao Wang, Julian Chan, Yuan Huang, Zhizheng Wu, Mingbo Ma</i>
RB-5 (#5)	SD-Eval: A Benchmark Dataset for Spoken Dialogue Understanding Beyond Words <i>Junyi Ao, Yuancheng Wang, Xiaohai Tian, Dekun Chen, Jun Zhang, Lu Lu, Yuxuan Wang, Haizhou Li, Zhizheng Wu</i>
RB-6 (#6)	SuperM2M: Supervised and Mixture-to-Mixture Co-Learning for Speech Enhancement and Robust ASR <i>Zhong-Qiu Wang</i>
RB-7 (#7)	Speech Recognition Corpus of the Khinalug Language for Documenting Endangered Languages <i>Zhaolin Li, Monika Rind-Pawłowski, Jan Niehues</i>
RB-8 (#8)	MaskGCT: Zero-Shot Text-to-Speech with Masked Generative Codec Transformer <i>Yuancheng Wang, Haoyue Zhan, Liwei Liu, Ruihong Zeng, Haotian Guo, Jiachen Zheng, Qiang Zhang, Xueyao Zhang, Shunsi Zhang, Zhizheng Wu</i>
RB-9 (#9)	Cacophony: An Improved Contrastive Audio-Text Model <i>Ge Zhu, Jordan Darefsky, Zhiyao Duan</i>
RB-10 (#10)	LlamaPartialSpoof: An LLM-Driven Fake Speech Dataset Simulating Disinformation Generation <i>Hieu-Thi Luong, Haoyang Li, Lin Zhang, Kong Aik Lee, Eng Siong Chng</i>
RB-11 (#11)	Self-Supervised ASR Models and Features for Dysarthric and Elderly Speech Recognition <i>Shujie Hu, Xurong Xie, Mengzhe Geng, Zengrui Jin, Jiajun Deng, Guinan Li, Yi Wang, Mingyu Cui, Tianzi Wang, Helen Meng, Xunying Liu</i>
RB-12 (#12)	PAT: Parameter-Free Audio-Text Aligner to Boost Zero-Shot Audio Classification <i>Ashish Seth, Ramaneswaran Selvakumar, Sonal Kumar, Sreyan Ghosh, Dinesh Manocha</i>
RB-13 (#13)	A Transformer Framework for Simultaneous Segmentation,

	<p>Classification, and Caller Identification of Marmoset Vocalization</p> <p><i>Bin Wu, Shinnosuke Takamichi, Sakriani Sakti, and Satoshi Nakamura</i></p>
RB-14 (#14)	<p>UTDUSS: UTokyo-SaruLab System for Interspeech2024 Speech Processing Using Discrete Speech Unit Challenge</p> <p><i>Wataru Nakata, Kazuki Yamauchi, Dong Yang, Hiroaki Hyodo, Yuki Saito</i></p>
RB-15 (#15)	<p>Optimizing Contextual Speech Recognition Using Vector Quantization for Efficient Retrieval</p> <p><i>Nikolaos Flemotomos, Roger Hsiao, Pawel Swietojanski, Takaaki Hori, Dogan Can, Xiaodan Zhuang</i></p>
RB-16 (#16)	<p>Streaming Speech-to-speech Speech-LLM for Simultaneous Translation and Multi-turn Conversation</p> <p><i>Elena Rastorgueva, Zhehuai Chen, He Huang, Edresson Casanova, Jason Li, Krishna Puvvada, Ryan Langman, Ante Jukić, Kunal Dhawan, Nithin Rao Koluguri, Subhankar Ghosh, Piotr Żelasko, Oleksii Hrinchuk, Andrei Andrusenko, Vitaly Lavrukhin, Jagadeesh Balam, Boris Ginsburg</i></p>

19:00-20:30 Welcome Reception

Day 2, Dec 3, Tuesday

09:00-10:00 Keynote Speech 2 (Venue: Lecture Hall)

Title: End-to-End Audio Processing: From On-Device Models to LLMs

Speaker: *Dr Tara N. Sainath, Google DeepMind*

Chair: *Prof Hung-yi Lee*

10:00-10:30 Coffee Break

10:30-12:30 Poster Session 3: Speech Processing (Venue: Poster Area)

Chair: Prof Junichi Yamagishi

Poster ID (Paper ID)	Title and Authors
P3-01-ANA (#66)	<p>Property Neurons in Self-Supervised Speech Transformers</p> <p><i>Tzu-Quan Lin (National Taiwan University)</i> <i>Guan-Ting Lin (National Taiwan University)</i> <i>Hung-Yi Lee (National Taiwan University)</i> <i>Hao Tang (The University of Edinburgh)</i></p>

<p>P3-02-ANA (#145)</p>	<p>Privacy vs Emotion Preservation Trade-Offs in Emotion-Preserving Speaker Anonymization</p> <p><i>Zexin Cai (Johns Hopkins University)</i> <i>Henry Li Xinyuan (Johns Hopkins University)</i> <i>Ashi Garg (Bharti Vidyapeeth College of Engineering)</i> <i>Nicholas O Andrews (Johns Hopkins University)</i> <i>Paola Garcia (Johns Hopkins University)</i> <i>Matthew S Wiesner (Johns Hopkins University)</i> <i>Kevin Duh (Johns Hopkins University)</i> <i>Sanjeev Khudanpur (Johns Hopkins University)</i></p>
<p>P3-03-ANA (#217)</p>	<p>Estimating the Completeness of Discrete Speech Units</p> <p><i>Sung-Lin Yeh (University of Edinburgh)</i> <i>Hao Tang (The University of Edinburgh)</i></p>
<p>P3-04-ANA (#314)</p>	<p>Investigation of Speaker Representation for Target-Speaker Speech Processing</p> <p><i>Takanori Ashihara (NTT Corp.)</i> <i>Takafumi Moriya (NTT Corporation)</i> <i>Shota Horiguchi (NTT Corporation)</i> <i>Junyi Peng (Brno University of Technology)</i> <i>Tsubasa Ochiai (NTT)</i> <i>Marc Delcroix (NTT)</i> <i>Kohei Matsuura (NTT)</i> <i>Hiroshi Sato (NTT Corporation)</i></p>
<p>P3-05-MMP (#226)</p>	<p>Crossmodal ASR Error Correction with Discrete Speech Units</p> <p><i>Yuanchao Li (University of Edinburgh)</i> <i>Pinzhen Chen (University of Edinburgh)</i> <i>Peter Bell (University of Edinburgh)</i> <i>Catherine Lai (University of Edinburgh)</i></p>
<p>P3-06-MMP (#241)</p>	<p>Listen and Speak Fairly: A Study on Semantic Gender Bias in Speech Integrated Large Language Models</p> <p><i>Yi-Cheng Lin (National Taiwan University)</i> <i>Tzu-Quan Lin (National Taiwan University)</i> <i>Chih-Kai Yang (National Taiwan University)</i> <i>Ke-Han Lu (National Taiwan University)</i> <i>Wei-Chih Chen (National Taiwan University)</i> <i>Chun-Yi Kuan (National Taiwan University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
<p>P3-07-MMP (#265)</p>	<p>Learning Video Temporal Dynamics with Cross-Modal Attention for Robust Audio-Visual Speech Recognition</p>

	<p><i>Sungnyun Kim (KAIST)</i> <i>Kangwook Jang (KAIST)</i> <i>Sangmin Bae (KAIST)</i> <i>Hoirin Kim (KAIST)</i> <i>Se-Young Yun (KAIST)</i></p>
P3-08-MMP (#269)	<p>Data Efficient Reflow for Few Step Audio Generation</p> <p><i>Lemeng Wu (Meta)</i> <i>Zhaoheng Ni (Meta AI)</i> <i>Bowen Shi (Toyota Technological Institute at Chicago)</i> <i>Wei-Ning Hsu (Meta)</i> <i>Gael Le Lan (Meta)</i> <i>Varun Nagaraja (Meta)</i> <i>Anurag Kumar (Meta)</i> <i>Xinhao Mei (Meta)</i> <i>Yunyang Xiong (Meta)</i> <i>Bilge Soran (Meta)</i> <i>Raghuraman Krishnamoorthi (Facebook)</i> <i>Yangyang Shi (Facebook)</i> <i>Vikas Chandra (Meta)</i></p>
P3-09-MLS (#99)	<p>Optimizing Byte-Level Representation for End-to-End ASR</p> <p><i>Roger Hsiao (Apple)</i> <i>Liuhui Deng (Apple)</i> <i>Erik McDermott (Apple)</i> <i>Ruchir Travadi (Apple)</i> <i>Xiaodan Zhuang (Apple)</i></p>
P3-10-MLS (#179)	<p>Romanization Encoding for Multilingual ASR</p> <p><i>Wen Ding (NVIDIA)</i> <i>Fei Jia (NVIDIA Corporation)</i> <i>Hainan Xu (NVIDIA Corporation)</i> <i>Yu Xi (NVIDIA)</i> <i>Junjie Lai (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i></p>
P3-11-MLS (#254)	<p>Enhancing Code-Switching ASR Leveraging Non-Peaky CTC Loss and Deep Language Posterior Injection</p> <p><i>Tzu-Ting Yang (National Taiwan Normal University)</i> <i>Hsin-Wei Wang (NTNU)</i> <i>Yi-Cheng Wang (National Taiwan Normal University)</i> <i>Berlin Chen (National Taiwan Normal University)</i></p>
P3-12-MLS	<p>Language-Independent Prosody-Enhanced Speech</p>

(#263)	<p>Representations for Multilingual Speech Synthesis</p> <p><i>Chang Liu (University of Science and Technology of China)</i> <i>Zhen-Hua Ling (University of Science and Technology of China)</i> <i>Ya-Jun Hu (iFLYTEK Co., Ltd.)</i></p>
P3-13-MLS (#359)	<p>Classification of Spontaneous and Scripted Speech for Multilingual Audio</p> <p><i>Shahar Elisha (Spotify)</i> <i>Andrew J McDowell (Spotify)</i> <i>Mariano Beguerisse-Díaz (Spotify)</i> <i>Emmanouil Benetos (Queen Mary University of London)</i></p>
P3-14-EMR (#40)	<p>GMP-TL: Gender-Augmented Multi-Scale Pseudo-Label Enhanced Transfer Learning for Speech Emotion Recognition</p> <p><i>Yu Pan (Kyushu University)</i> <i>Yuguang Yang (Ximalaya Inc., Shanghai, China)</i> <i>Yuheng Huang (The University of Tokyo)</i> <i>Tiancheng Jin (Kyushu University)</i> <i>Jingjing Yin (Ximalaya)</i> <i>Yanni Hu (Ximalaya Inc., Shanghai, China)</i> <i>Heng Lu (Ximalaya Inc.)</i> <i>Lei Ma (The University of Tokyo / University of Alberta)</i> <i>Jianjun Zhao (Kyushu University)</i></p>
P3-15-EMR (#81)	<p>Embracing Ambiguity and Subjectivity Using the All-Inclusive Aggregation Rule for Evaluating Multi-Label Speech Emotion Recognition Systems</p> <p><i>Huang-Cheng Chou (Department of Electrical Engineering at National Tsing Hua University (NTHU))</i> <i>Haibin Wu (National Taiwan University)</i> <i>Lucas Goncalves (The University of Texas at Dallas)</i> <i>Seong-Gyun Leem (University of Texas at Dallas)</i> <i>Ali N Salman (University of Texas at Dallas)</i> <i>Carlos Busso (University of Texas at Dallas)</i> <i>Hung-Yi Lee (National Taiwan University)</i> <i>Chi-Chun Lee (National Tsing Hua University)</i></p>
P3-16-EMR (#83)	<p>Open-Emotion: A Reproducible EMO-SUPERB for Speech Emotion Recognition Systems</p> <p><i>Haibin Wu (National Taiwan University)</i> <i>Huang-Cheng Chou (Department of Electrical Engineering at National Tsing Hua University (NTHU))</i> <i>Kai-Wei Chang (National Taiwan University)</i> <i>Lucas Goncalves (The University of Texas at Dallas)</i> <i>Jiawei Du (National Taiwan University)</i></p>

	<p><i>Jyh-Shing Roger Jang (National Taiwan University)</i> <i>Chi-Chun Lee (National Tsing Hua University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P3-17-EMR (#225)	<p>Speech Emotion Recognition with ASR Transcripts: A Comprehensive Study on Word Error Rate and Fusion Techniques</p> <p><i>Yuanchao Li (University of Edinburgh)</i> <i>Peter Bell (University of Edinburgh)</i> <i>Catherine Lai (University of Edinburgh)</i></p>
P3-18-EMR (#282)	<p>Beyond the Binary: Limitations and Possibilities of Gender-Related Speech Technology Research</p> <p><i>Ariadna Sanchez (The University of Edinburgh)</i> <i>Alice Ross (University of Edinburgh)</i> <i>Nina Markl (University of Essex)</i></p>
P3-19-EMR (#352)	<p>Enhancing Domain Generalization in Speech Emotion Recognition by Combining Domain-Variant Representations and Domain-Invariant Classifiers</p> <p><i>Shi-Wook Lee (National Institute of Advanced Industrial Science and Technology)</i></p>
P3-20-SS07 (#47)	<p>MDCTCodec: A Lightweight MDCT-Based Neural Audio Codec for High Sampling Rate and Low Bitrate Scenarios</p> <p><i>Xiao-Hang Jiang (University of Science and Technology of China)</i> <i>Yang Ai (University of Science and Technology of China)</i> <i>Rui-Chen Zheng (University of Science and Technology of China)</i> <i>Hui-Peng Du (University of Science and Technology of China)</i> <i>Ye-Xin Lu (University of Science and Technology of China)</i> <i>Zhen-Hua Ling (University of Science and Technology of China)</i></p>
P3-21-SS07 (#51)	<p>Addressing Index Collapse of Large-Codebook Speech Tokenizer with Dual-Decoding Product-Quantized Variational Auto-Encoder</p> <p><i>Haohan Guo (The Chinese University of Hong Kong)</i> <i>Fenglong Xie (Xiaohongshu)</i> <i>Dongchao Yang (The Chinese University of Hong Kong)</i> <i>Hui Lu (The Chinese University of Hong Kong)</i> <i>Xixin Wu (The Chinese University of Hong Kong)</i> <i>Helen Meng (The Chinese University of Hong Kong)</i></p>
P3-22-SS07 (#267)	<p>Investigating Neural Audio Codecs for Speech Language Model-Based Speech Generation</p> <p><i>Jiaqi Li (The Chinese University of Hong Kong, Shenzhen)</i></p>

	<p> <i>Dongmei Wang (Microsoft)</i> <i>Xiaofei Wang (Microsoft)</i> <i>Yao Qian (Microsoft)</i> <i>Long Zhou (Microsoft Research Asia)</i> <i>Shujie Liu (Microsoft Research Asia)</i> <i>Midia Yousefi (Microsoft)</i> <i>Canrun Li (Microsoft)</i> <i>Chung-Hsien Tsai (Microsoft)</i> <i>Zhen Xiao (Microsoft)</i> <i>Yanqing Liu (Microsoft)</i> <i>Junkun Chen (Microsoft)</i> <i>Sheng Zhao (Microsoft)</i> <i>Jinyu Li (Microsoft)</i> <i>Zhizheng Wu (Chinese University of Hong Kong, Shenzhen)</i> <i>Michael Zeng (Microsoft)</i> </p>
<p>P3-23-SS07 (#280)</p>	<p>ESPnet-Codex: Comprehensive Training and Evaluation of Neural Codexes for Audio, Music, and Speech</p> <p> <i>Jiatong Shi (Carnegie Mellon University)</i> <i>Jinchuan Tian (Carnegie Mellon University)</i> <i>Yihan Wu (Renmin University of China)</i> <i>Jee-Weon Jung (Carnegie Mellon University)</i> <i>Jia Qi Yip (Alibaba Group / Nanyang Technological University)</i> <i>Yoshiki Masuyama (Tokyo Metropolitan University)</i> <i>William Chen (Carnegie Mellon University)</i> <i>Yuning Wu (Renmin University of China)</i> <i>Yuxun Tang (Renmin University of China)</i> <i>Massa Baali (CMU)</i> <i>Dareen Alharthi (Carnegie Mellon University)</i> <i>Dong Zhang (Fudan University)</i> <i>Ruifan Deng (Fudan University)</i> <i>Tejes Srivastava (University of Chicago)</i> <i>Haibin Wu (National Taiwan University)</i> <i>Alexander H Liu (MIT)</i> <i>Bhiksha Raj (Carnegie Mellon University)</i> <i>Qin Jin (Renmin University of China)</i> <i>Ruihua Song (Renmin University of China)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> </p>
<p>P3-24-SS07 (#336)</p>	<p>Codex-SUPERB @ SLT 2024: A Lightweight Benchmark for Neural Codex Models</p> <p> <i>Haibin Wu (National Taiwan University)</i> <i>Xuanjun Chen (National Taiwan University)</i> <i>Yi-Cheng Lin (National Taiwan University)</i> <i>Jiawei Du (National Taiwan University)</i> <i>Kai-Wei Chang (National Taiwan University)</i> <i>Ke-Han Lu (National Taiwan University)</i> <i>Alexander H Liu (MIT)</i> </p>

	<p><i>Ho Lam Chung (National Taiwan University)</i> <i>Yuan-Kuei Wu (National Taiwan University)</i> <i>Dongchao Yang (The Chinese University of Hong Kong)</i> <i>Songxiang Liu (Tencent)</i> <i>Yi-Chiao Wu (Meta)</i> <i>Xu Tan (Microsoft Research Asia)</i> <i>James Glass (Massachusetts Institute of Technology)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P3-25-SS08 (#61)	<p>Optimizing Dysarthria Wake-Up Word Spotting: An End-to-End Approach for SLT 2024 LRDWWS Challenge</p> <p><i>Shuiyun Liu (Northwestern Polytechnical University)</i> <i>Yuxiang Kong (Xiaomi Inc.)</i> <i>Pengcheng Guo (Northwestern Polytechnical University)</i> <i>Weiji Zhuang (Xiaomi Inc.)</i> <i>Peng Gao (Xiaomi Inc.)</i> <i>Yujun Wang (Xiaomi)</i> <i>Lei Xie (NWPU)</i></p>
P3-26-SS08 (#234)	<p>PB-LRDWWS System for the SLT 2024 Low-Resource Dysarthria Wake-Up Word Spotting Challenge</p> <p><i>Shiyao Wang (Nankai University)</i> <i>Jiaming Zhou (Nankai University)</i> <i>Shiwan Zhao (Nankai University)</i> <i>Yong Qin (Nankai University)</i></p>
P3-27-SS08 (#404)	<p>Summary of Low-Resource Dysarthria Wake-Up Word Spotting Challenge</p> <p><i>Ming Gao (University of Science and Technology of China)</i> <i>Hang Chen (USTC)</i> <i>Jun Du (University of Science and Technology of China)</i> <i>Xin Xu (Beijing AISHELL Technology Co., Ltd.)</i> <i>Hongxiao Guo (Beijing AISHELL Technology Co., Ltd.)</i> <i>Hui Bu (AISHELL)</i> <i>Ming Li (Wuhan University)</i> <i>Chin-Hui Lee (Georgia Institute of Technology)</i></p>
P3-28-SLP (#105)	<p>ProGRES: Prompted Generative Rescoring on ASR N-Best</p> <p><i>Ada D Tur (McGill University)</i> <i>Mirco Ravanelli (Concordia University, Université de Montréal, MILA)</i> <i>Adel Moumen (Avignon University)</i></p>
P3-29-SS02 (#212)	<p>FLANEC: Exploring Flan-T5 for Post-ASR Error Correction</p>

	<p>Moreno La Quatra (<i>Kore University of Enna</i>) Valerio Mario Salerno (<i>Università degli Studi di Enna "Kore"</i>) Yu Tsao (<i>Accademia Sinica</i>) Sabato Marco Siniscalchi (<i>Kore University of Enna</i>)</p>
--	--

**10:30-12:30 Challenge Session 3: Source Speaker Tracing Challenge(SSTC)
(Venue: Lecture Hall)**

12:30-14:00 Lunch

14:00-15:00 Invited Talk 2 (Venue: Lecture Hall)

Title: Holistic Artificial Intelligence (HAI): From Big Models to Big Applications

Speaker: *Dr Junlan Feng, China Mobile Research Institute*

Chair: *Prof Lei Wang*

15:00-15:30 Coffee Break

15:30-17:30 Poster Session 4: Speech Synthesis (Venue: Poster Area)

Chair: Prof Xixin Wu

Poster ID (Paper ID)	Title and Authors
P4-01-TTS (#21)	<p>AS-Speech: Adaptive Style for Speech Synthesis</p> <p><i>Zhipeng Li (South China University of Technology)</i> <i>Xiaofen Xing (South China University of Technology)</i> <i>Jun Wang (Meituan)</i> <i>Shuaiqi Chen (South China University of Technology)</i> <i>Guoqiao Yu (Meituan)</i> <i>Guanglu Wan (Meituan)</i> <i>Xiangmin Xu (South China University of Technology)</i></p>
P4-02-TTS (#31)	<p>Room Impulse Responses Help Attackers Evade Deep Fake Detection</p> <p><i>Hieu-Thi Luong (Nanyang Technological University)</i> <i>Duc-Tuan Truong (Nanyang Technological University)</i> <i>Kong Aik Lee (The Hong Kong Polytechnic University)</i> <i>Eng Siong Chng (Nanyang Technological University)</i></p>
P4-03-TTS (#35)	<p>Attention-Constrained Inference for Robust Decoder-Only Text-to-Speech</p> <p><i>Hankun Wang (Shanghai Jiao Tong University)</i> <i>Chenpeng Du (Shanghai Jiao Tong University)</i> <i>Yiwei Guo (Shanghai Jiao Tong University)</i> <i>Shuai Wang (Shenzhen Research Institute of Big Data, Chinese University)</i></p>

	<p><i>of Hong Kong (Shenzhen)</i> <i>Xie Chen (Shanghai Jiao Tong University)</i> <i>Kai Yu (Shanghai Jiao Tong University)</i></p>
P4-04-TTS (#46)	<p>Stage-Wise and Prior-Aware Neural Speech Phase Prediction</p> <p><i>Fei Liu (University of Science and Technology of China)</i> <i>Yang Ai (University of Science and Technology of China)</i> <i>Hui-Peng Du (University of Science and Technology of China)</i> <i>Ye-Xin Lu (University of Science and Technology of China)</i> <i>Rui-Chen Zheng (University of Science and Technology of China)</i> <i>Zhen-Hua Ling (University of Science and Technology of China)</i></p>
P4-05-TTS (#52)	<p>SoCodec: A Semantic-Ordered Multi-Stream Speech Codec for Efficient Language Model-Based Text-to-Speech Synthesis</p> <p><i>Haohan Guo (The Chinese University of Hong Kong)</i> <i>Fenglong Xie (Xiaohongshu)</i> <i>Dongchao Yang (The Chinese University of Hong Kong)</i> <i>Xixin Wu (The Chinese University of Hong Kong)</i> <i>Helen Meng (The Chinese University of Hong Kong)</i> <i>Kun Xie (Xiaohongshu)</i> <i>Dake Guo (Northwestern Polytechnical University)</i></p>
P4-06-TTS (#56)	<p>Detecting the Undetectable: Assessing the Efficacy of Current Spoof Detection Methods Against Seamless Speech Edits</p> <p><i>Sung-Feng Huang (National Taiwan University)</i> <i>Heng-Cheng Kuo (National Taiwan University)</i> <i>Zhehuai Chen (NVIDIA)</i> <i>Xuesong Yang (NVIDIA Applied AI Research)</i> <i>Chao-Han Huck Yang (NVIDIA Research)</i> <i>Yu Tsao (Academia Sinica)</i> <i>Yu-Chiang Frank Wang (National Taiwan University)</i> <i>Hung-Yi Lee (National Taiwan University)</i> <i>Szu-Wei Fu (NVIDIA)</i></p>
P4-07-TTS (#62)	<p>DNN-Based Ensemble Singing Voice Synthesis with Interactions Between Singers</p> <p><i>Hiroaki Hyodo (The University of Tokyo)</i> <i>Shinnosuke Takamichi (Keio University)</i> <i>Tomohiko Nakamura (National Institute of Advanced Industrial Science and Technology (AIST))</i> <i>Junya Koguchi (Meiji University)</i> <i>Hiroshi Saruwatari (The University of Tokyo)</i></p>
P4-08-TTS (#87)	<p>Investigating Disentanglement in a Phoneme-Level Speech Codec for Prosody Modeling</p>

	<p><i>Sotirios Karapiperis (Samsung)</i> <i>Nikolaos Ellinas (Innoetics, Samsung Electronics)</i> <i>Alexandra Vioni (Innoetics, Samsung Electronics)</i> <i>Junkwang Oh (Mobile eXperience Business, Samsung Electronics)</i> <i>Gunu Jho (Mobile eXperience Business, Samsung Electronics)</i> <i>Inchul Hwang (Samsung Research)</i> <i>Spyros Raptis (Samsung Electronics Hellas / INNOETICS)</i></p>
P4-09-TTS (#128)	<p>InstructSing: High-Fidelity Singing Voice Generation via Instructing Yourself</p> <p><i>Chang Zeng (National Institute of Informatics)</i> <i>Chunhui Wang (Geely Automobile Research Institute)</i> <i>Xiaoxiao Miao (Singapore Institute of Technology)</i> <i>Jian Zhao (Geely)</i> <i>Zhonglin Jiang (Geely)</i> <i>Yong Chen (Geely Automobile Research Institute)</i></p>
P4-10-TTS (#166)	<p>E2 TTS: Embarrassingly Easy Fully Non-Autoregressive Zero-Shot TTS</p> <p><i>Sefik Emre Eskimez (Microsoft)</i> <i>Xiaofei Wang (Microsoft)</i> <i>Manthan Thakker (Microsoft)</i> <i>Canrun Li (Microsoft)</i> <i>Chung-Hsien Tsai (Microsoft)</i> <i>Zhen Xiao (Microsoft)</i> <i>Hemin Yang (Microsoft)</i> <i>Zirun Zhu (Microsoft)</i> <i>Min Tang (Microsoft)</i> <i>Xu Tan (Microsoft Research Asia)</i> <i>Yanqing Liu (Microsoft)</i> <i>Sheng Zhao (Microsoft)</i> <i>Naoyuki Kanda (Microsoft)</i></p>
P4-11-TTS (#167)	<p>Laugh Now Cry Later: Controlling Time-Varying Emotional States of Flow-Matching-Based Zero-Shot Text-to-Speech</p> <p><i>Haibin Wu (National Taiwan University)</i> <i>Xiaofei Wang (Microsoft)</i> <i>Sefik Emre Eskimez (Microsoft)</i> <i>Manthan Thakker (Microsoft)</i> <i>Daniel Tompkins (Microsoft)</i> <i>Chung-Hsien Tsai (Microsoft)</i> <i>Canrun Li (Microsoft)</i> <i>Zhen Xiao (Microsoft)</i> <i>Sheng Zhao (Microsoft)</i> <i>Jinyu Li (Microsoft)</i> <i>Naoyuki Kanda (Microsoft)</i></p>

P4-12-TTS (#184)	<p>Disentangling the Prosody and Semantic Information with Pre-Trained Model for In-Context Learning-Based Zero-Shot Voice Conversion</p> <p><i>Zhengyang Chen (Shanghai Jiao Tong University)</i> <i>Shuai Wang (Shenzhen Research Institute of Big Data, Chinese University of Hong Kong (Shenzhen))</i> <i>Mingyang Zhang (Chinese University of Hong Kong, Shenzhen)</i> <i>Xuechen Liu (National Institute of Informatics)</i> <i>Junichi Yamagishi (National Institute of Informatics)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i></p>
P4-13-TTS (#195)	<p>NDVQ: Robust Neural Audio Codec with Distribution-Based Vector Quantization</p> <p><i>Zhikang Niu (Shanghai Jiao Tong University)</i> <i>Sanyuan Chen (Harbin Institute of Technology)</i> <i>Long Zhou (Microsoft Research Asia)</i> <i>Ziyang Ma (Shanghai Jiao Tong University)</i> <i>Xie Chen (Shanghai Jiao Tong University)</i> <i>Shujie Liu (Microsoft Research Asia)</i></p>
P4-14-TTS (#228)	<p>Fast, High-Quality, and Parameter-Efficient Articulatory Synthesis Using Differentiable DSP</p> <p><i>Yisi Liu (University of California, Berkeley)</i> <i>Bohan Yu (UC Berkeley)</i> <i>Drake Lin (UC Berkeley)</i> <i>Peter Wu (UC Berkeley)</i> <i>Cheol Jun Cho (UC Berkeley)</i> <i>Gopala Krishna Anumanchipalli (UC Berkeley)</i></p>
P4-15-TTS (#299)	<p>VISinger2+: End-to-End Singing Voice Synthesis Augmented by Self-Supervised Learning Representation</p> <p><i>Yifeng Yu (Georgia Institute of Technology)</i> <i>Jiatong Shi (Carnegie Mellon University)</i> <i>Yuning Wu (Renmin University of China)</i> <i>Yuxun Tang (Renmin University of China)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i></p>
P4-16-TTS (#316)	<p>End-to-End Streaming Model for Low-Latency Speech Anonymization</p> <p><i>Waris Quamer (Texas A&M University)</i> <i>Ricardo Gutierrez-Osuna (Texas A&M University)</i></p>
P4-17-TTS	<p>Emotion-Coherent Speech Data Augmentation and Self-Supervised</p>

(#326)	<p>Contrastive Style Training for Enhancing Kids' Story Speech Synthesis</p> <p><i>Raymond Chung (LSCM)</i></p>
P4-18-TTS (#332)	<p>Discrete Unit-Based Masking for Improving Disentanglement in Voice Conversion</p> <p><i>Philip Lee (University of Texas at Dallas)</i> <i>İsmail Rasim Ülgen (University of Texas at Dallas)</i> <i>Berrak Sisman (University of Texas at Dallas)</i></p>
P4-19-TTS (#345)	<p>Cross-Dialect Text-To-Speech in Pitch-Accent Language Incorporating Multi-Dialect Phoneme-Level BERT</p> <p><i>Kazuki Yamauchi (The University of Tokyo)</i> <i>Yuki Saito (The University of Tokyo, Japan)</i> <i>Hiroshi Saruwatari (The University of Tokyo)</i></p>
P4-20-TTS (#353)	<p>Leveraging Diverse Semantic-Based Audio Pretrained Models for Singing Voice Conversion</p> <p><i>Xueyao Zhang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Zihao Fang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Yicheng Gu (The Chinese University of Hong Kong, Shenzhen)</i> <i>Haopeng Chen (The Chinese University of Hong Kong, Shenzhen)</i> <i>Lexiao Zou (Harbin Institute of Technology (Shenzhen))</i> <i>Junan Zhang (Fudan University)</i> <i>Liumeng Xue (The Chinese University of Hong Kong, Shenzhen)</i> <i>Zhizheng Wu (The Chinese University of Hong Kong, Shenzhen)</i></p>
P4-21-TTS (#394)	<p>TTSDS: Text-to-Speech Distribution Score</p> <p><i>Christoph D Minixhofer (The University of Edinburgh)</i> <i>Ondřej Klejch (University of Edinburgh)</i> <i>Peter Bell (University of Edinburgh)</i></p>
P4-22-SS03 (#261)	<p>Speech Foundation Model Ensembles for the Controlled Singing Voice Deepfake Detection (CtrSVDD)</p> <p><i>Anmol Guragain (Vellore Institute of Technology)</i> <i>Tianchi Liu (National University of Singapore)</i> <i>Zihan Pan (Institute for Infocomm Research (I2R), ASTAR, Singapore)</i> <i>Hardik B Sailor (I2R, ASTAR, Singapore)</i> <i>Qiongqiong Wang (ASTAR)</i></p>
P4-23-SS03 (#323)	<p>SVDD 2024: The Inaugural Singing Voice Deepfake Detection Challenge</p>

	<p><i>You Zhang (University of Rochester)</i> <i>Yongyi Zang (University of Rochester)</i> <i>Jiatong Shi (Carnegie Mellon University)</i> <i>Ryuichi Yamamoto (Nagoya University)</i> <i>Tomoki Toda (Nagoya University)</i> <i>Zhiyao Duan (University of Rochester)</i></p>
P4-24-SS03 (#348)	<p>XWSB: A Blend System Utilizing XLS-R and WavLM with SLS Classifier for the SVDD 2024 Challenge</p> <p><i>Zhang Qishan (Hubei Minzu University)</i> <i>Shuangbing Wen (Hubei Minzu University)</i> <i>Fangke Yan (Hubei Minzu University)</i> <i>Tao Hu (Hubei Minzu University)</i> <i>Jun Li (Hubei Minzu University)</i></p>
P4-25-SS03 (#416)	<p>Integrating Self-Supervised Pre-Training with Adversarial Learning for Synthesized Song Detection</p> <p><i>Yankai Wang (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Yuxuan Du (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Dejun Zhang (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Rong Zheng (Beijing Fosafer Information Technology Co., Ltd.)</i> <i>Jing Deng (Beijing Fosafer Information Technology Co., Ltd.)</i></p>
P4-26-SS05 (#396)	<p>The VoiceMOS Challenge 2024: Beyond Speech Quality Prediction</p> <p><i>Wen-Chin Huang (Nagoya University)</i> <i>Szu-Wei Fu (NVIDIA)</i> <i>Erica Cooper (National Institute of Information and Communications Technology)</i> <i>Ryandhimas E Zezario (Academia Sinica)</i> <i>Tomoki Toda (Nagoya University)</i> <i>Hsin-Min Wang (Academia Sinica)</i> <i>Junichi Yamagishi (National Institute of Informatics)</i> <i>Yu Tsao (Academia Sinica)</i></p>
P4-27-SS05 (#406)	<p>Pitch-and-Spectrum-Aware Singing Quality Assessment with Bias Correction and Model Fusion</p> <p><i>Yu-Fei Shi (University of Science and Technology of China)</i> <i>Yang Ai (University of Science and Technology of China)</i> <i>Ye-Xin Lu (University of Science and Technology of China)</i> <i>Hui-Peng Du (University of Science and Technology of China)</i> <i>Zhen-Hua Ling (University of Science and Technology of China)</i></p>
P4-28-SS05 (#407)	<p>The T05 System for The VoiceMOS Challenge 2024: Transfer Learning from Deep Image Classifier to Naturalness MOS Prediction of High-Quality Synthetic Speech</p>

	<p><i>Kaito Baba (The University of Tokyo)</i> <i>Wataru Nakata (The University of Tokyo)</i> <i>Yuki Saito (The University of Tokyo, Japan)</i> <i>Hiroshi Saruwatari (The University of Tokyo)</i></p>
--	---

15:15-17:30 Challenge Session 4: Codec SUPERB Challenge (Venue: Lecture Hall)

Day 3, Dec 4, Wednesday

09:00-10:00 Keynote Speech 3 (Venue: Lecture Hall)

Title: Large Language-Audio Models and Applications

Speaker: Prof Wenwu Wang, University of Surrey

Chair: Dr Jinyu Li

10:00-10:30 Coffee Break

10:30-12:30 Poster Session 5: Machine Learning & Resources (Venue: Poster Area)

Chair: Prof Ming Li

Poster ID (Paper ID)	Title and Authors
P5-01-TLP (#12)	<p>Automated Speaking Assessment of Conversation Tests with a Novel Graph-Based Modeling Method on Spoken Response Coherence</p> <p><i>Jiun-Ting Li (National Taiwan Normal University)</i> <i>Bi-Cheng Yan (National Taiwan Normal University)</i> <i>Tien-Hong Lo (National Taiwan Normal University)</i> <i>Yi-Cheng Wang (National Taiwan Normal University)</i> <i>Yung-Chang Hsu (EZ-AI)</i> <i>Berlin Chen (National Taiwan Normal University)</i></p>
P5-02-TLP (#115)	<p>Conditional Label Smoothing for LLM-Based Data Augmentation in Medical Text Classification</p> <p><i>Luca Manuel Becker (Institute of Communication Acoustics, Ruhr-Universität Bochum)</i> <i>Philip Pracht (Bochum Institute of Technology)</i> <i>Peter Sertdal (Fraunhofer Institute for High Frequency Physics and Radar Techniques FHR)</i> <i>Jil Uboreck (Bochum Institute of Technology)</i> <i>Alexander Bendel (Institute of Work and Qualification, University of</i></p>

	<p><i>Duisburg-Essen</i> <i>Rainer Martin (Institute of Communication Acoustics, Ruhr-Universität Bochum)</i></p>
P5-03-TLP (#180)	<p>Plan, Generate and Optimize: Extending Large Language Models for Dialogue Systems via Prompt-Based Collaborative Method</p> <p><i>Mengfei Guo (China Mobile Research Institute (CMRI))</i> <i>Si Chen (China Mobile Research Institute (CMRI))</i> <i>Yi Huang (China Mobile Research)</i> <i>Junlan Feng (China Mobile Research)</i></p>
P5-04-TLP (#202)	<p>Taming NLU Noise: Student-Teacher Learning for Robust Dialogue Policy</p> <p><i>Mahdin Rohmatillah (Universitas Brawijaya)</i> <i>Jen-Tzung Chien (National Yang Ming Chiao Tung University)</i></p>
P5-05-RES (#15)	<p>HeightCeleb: An Enrichment of VoxCeleb Dataset with Speaker Height Information</p> <p><i>Stanisław Kacprzak (AGH University of Krakow)</i> <i>Konrad Kowalczyk (AGH University of Krakow)</i></p>
P5-06-RES (#57)	<p>ESPnet-EZ: Python-Only ESPnet for Easy Fine-Tuning and Integration</p> <p><i>Masao Someki (IBM)</i> <i>Kwanghee Choi (Carnegie Mellon University)</i> <i>Siddhant Arora (Carnegie Mellon University)</i> <i>William Chen (Carnegie Mellon University)</i> <i>Samuele Cornell (Carnegie Mellon University)</i> <i>Jionghao Han (Carnegie Mellon University)</i> <i>Yifan Peng (Carnegie Mellon University)</i> <i>Jiatong Shi (Carnegie Mellon University)</i> <i>Vaibhav Srivastav (Hugging Face, Inc.)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i></p>
P5-07-RES (#106)	<p>Spoken Stereaset: On Evaluating Social Bias Toward Speaker in Speech Large Language Models</p> <p><i>Yi-Cheng Lin (National Taiwan University)</i> <i>Wei-Chih Chen (National Taiwan University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P5-08-RES (#124)	<p>Amphion: An Open-Source Audio, Music, and Speech Generation Toolkit</p> <p><i>Xueyao Zhang (The Chinese University of Hong Kong, Shenzhen)</i></p>

	<p> <i>Liუმeng Xue (The Chinese University of Hong Kong, Shenzhen)</i> <i>Yicheng Gu (The Chinese University of Hong Kong, Shenzhen)</i> <i>Yuancheng Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Jiaqi Li (The Chinese University of Hong Kong, Shenzhen)</i> <i>Haorui He (The Chinese University of Hong Kong, Shenzhen)</i> <i>Chaoren Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Liu Songting (Nanyang Technological University)</i> <i>Xi Chen (Chinese University of Hong Kong (Shenzhen))</i> <i>Junan Zhang (Fudan University)</i> <i>Zihao Fang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Haopeng Chen (The Chinese University of Hong Kong, Shenzhen)</i> <i>Tze Ying Tang (CUHK-Shenzhen)</i> <i>Lexiao Zou (Harbin Institute of Technology (Shenzhen))</i> <i>Mingxuan Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Jun Han (The Chinese University of Hong Kong, Shenzhen)</i> <i>Kai Chen (Shanghai AI Laboratory)</i> <i>Haizhou Li (The Chinese University of Hong Kong, Shenzhen)</i> <i>Zhizheng Wu (The Chinese University of Hong Kong, Shenzhen)</i> </p>
P5-09-RES (#126)	<p> Emilia: An Extensive, Multilingual, and Diverse Speech Dataset for Large-Scale Speech Generation </p> <p> <i>Haorui He (The Chinese University of Hong Kong, Shenzhen)</i> <i>Zengqiang Shang (The Institute of Acoustics of the Chinese Academy of Sciences)</i> <i>Chaoren Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Xuyuan Li (The Institute of Acoustics of the Chinese Academy of Sciences)</i> <i>Yicheng Gu (The Chinese University of Hong Kong, Shenzhen)</i> <i>Hua Hua (Key Laboratory of Speech Acoustics and Content Understanding, Institute of Acoustics, China)</i> <i>Liwei Liu (The Chinese University of Hong Kong, Shenzhen)</i> <i>Chen Yang (The Institute of Acoustics of the Chinese Academy of Sciences)</i> <i>Jiaqi Li (The Chinese University of Hong Kong, Shenzhen)</i> <i>Peiyang Shi (The Institute of Acoustics of the Chinese Academy of Sciences)</i> <i>Yuancheng Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Kai Chen (Shanghai AI Laboratory)</i> <i>Pengyuan Zhang (Institute of Acoustics, Chinese Academy of Sciences)</i> <i>Zhizheng Wu (The Chinese University of Hong Kong, Shenzhen)</i> </p>
P5-10-RES (#191)	<p> FLORAS 50: A Massively Multilingual Multitask Benchmark for Long-Form Conversational Speech </p> <p> <i>William Chen (Carnegie Mellon University)</i> <i>Brian Yan (Carnegie Mellon University)</i> <i>Chih-Chen Chen (TMU)</i> <i>Shinji Watanabe (Carnegie Mellon University)</i> </p>
P5-11-RES	<p> Massively Multilingual Forced Aligner Leveraging Self-Supervised </p>

<p>(#295)</p>	<p>Discrete Units</p> <p><i>Hirofumi Inaguma (Meta)</i> <i>Iliia Kulikov (Meta)</i> <i>Zhaoheng Ni (Meta AI)</i> <i>Sravya Popuri (Meta)</i> <i>Paden P Tomasello (Meta)</i></p>
<p>P5-12-RES (#298)</p>	<p>Speech Recognition for Analysis of Police Radio Communication</p> <p><i>Tejes Srivastava (University of Chicago)</i> <i>Ju-Chieh Chou (TTIC)</i> <i>Priyank Shroff (University of Chicago)</i> <i>Karen Livescu (TTI-Chicago)</i> <i>Christopher Graziul (University of Chicago)</i></p>
<p>P5-13-RES (#307)</p>	<p>Large Language Models as User-Agents for Evaluating Task-Oriented Dialogue Systems</p> <p><i>Taaha Kazi (University of Illinois at Urbana-Champaign)</i> <i>Ruilang Lyu (University of Illinois at Urbana-Champaign)</i> <i>Sizhe Zhou (University of Illinois at Urbana-Champaign)</i> <i>Dilek Hakkani-Tur (University of Illinois, Urbana-Champaign)</i> <i>Gokhan Tur (Amazon)</i></p>
<p>P5-14-RES (#334)</p>	<p>DFADD: The Diffusion and Flow-Matching Based Audio Deepfake Dataset</p> <p><i>Jiawei Du (National Taiwan University)</i> <i>I-Ming Lin (National Taiwan University)</i> <i>I-Hsiang Chiu (National Taiwan University)</i> <i>Xuanjun Chen (National Taiwan University)</i> <i>Haibin Wu (National Taiwan University)</i> <i>Wenze Ren (National Taiwan University)</i> <i>Yu Tsao (Academia Sinica)</i> <i>Hung-Yi Lee (National Taiwan University)</i> <i>Roger Jang</i></p>
<p>P5-15-RES (#384)</p>	<p>SpMis: An Investigation of Synthetic Spoken Misinformation Detection</p> <p><i>Peizhuo Liu (The Chinese University of Hong Kong, Shenzhen)</i> <i>Li Wang (The Chinese University of Hong Kong, Shenzhen)</i> <i>He Renqiang (The Chinese University of Hong Kong, Shenzhen)</i> <i>Haorui He (The Chinese University of Hong Kong, Shenzhen)</i> <i>Lei Wang (Huawei International)</i> <i>Huadi Zheng (Huawei Technology)</i> <i>Jie Shi (Huawei International)</i> <i>Tong Xiao (Northeastern University)</i> <i>Zhizheng Wu (The Chinese University of Hong Kong, Shenzhen)</i></p>

P5-16-MLS (#7)	<p>Self-Supervised Speech Models for Word-Level Stuttered Speech Detection</p> <p><i>Yi-Jen Shih (The University of Texas at Austin)</i> <i>Zoi Gkalitsiou (UT Austin)</i> <i>Alex Dimakis (UT Austin)</i> <i>David Harwath (The University of Texas at Austin)</i></p>
P5-17-MLS (#30)	<p>Enhancing Automatic Speech Assessment Leveraging Heterogeneous Features and Soft Labels for Ordinal Classification</p> <p><i>Wen Hsuan Peng (National Taiwan Normal University)</i> <i>Sally Chen (The Language Training & Testing Center)</i> <i>Berlin Chen (National Taiwan Normal University)</i></p>
P5-18-MLS (#92)	<p>Speech Recognition-Based Feature Extraction for Enhanced Automatic Severity Classification in Dysarthric Speech</p> <p><i>Jeehyun Lee (Sogang University)</i> <i>Yerin Choi (Sogang University)</i> <i>Myoung-Wan Koo (Sogang University)</i></p>
P5-19-MLS (#144)	<p>Efficient Training of Self-Supervised Speech Foundation Models on a Compute Budget</p> <p><i>Andy T. Liu (National Taiwan University)</i> <i>Yi-Cheng Lin (National Taiwan University)</i> <i>Haibin Wu (National Taiwan University)</i> <i>Stefan Winkler (National University of Singapore)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P5-20-MLS (#175)	<p>Improving Anomalous Sound Detection via Low-Rank Adaptation Fine-Tuning of Pre-Trained Audio Models</p> <p><i>Xinhu Zheng (Tsinghua University)</i> <i>Anbai Jiang (Tsinghua University)</i> <i>Bing Han (Shanghai Jiao Tong University)</i> <i>Yanmin Qian (Shanghai Jiao Tong University)</i> <i>Pingyi Fan (Tsinghua University)</i> <i>Jia Liu (Tsinghua University)</i> <i>Wei-Qiang Zhang (Tsinghua University)</i></p>
P5-21-MLS (#233)	<p>Exploring ASR-Based Wav2Vec2 for Automated Speech Disorder Assessment: Insights and Analysis</p> <p><i>Tuan Manh Nguyen (LIA, Avignon University)</i> <i>Corinne Fredouille (Avignon Université- LIA)</i> <i>Alain Ghio (Aix-Marseille University, LPL)</i></p>

	<p><i>Mathieu Balaguer (IRIT)</i> <i>Virginie Woisard (Hospitals of Toulouse)</i></p>
P5-22-MLS (#249)	<p>Hierarchical Multi-Path and Multi-Model Selection for Fake Speech Detection</p> <p><i>Chang Feng (Tsinghua University)</i> <i>Yiyang Zhao (Tsinghua University)</i> <i>Guangzhi Sun (University of Cambridge Department of Engineering)</i> <i>Zehua Chen (Tsinghua University)</i> <i>Shuai Wang (Shenzhen Research Institute of Big Data, Chinese University of Hong Kong (Shenzhen))</i> <i>Chao Zhang (Tsinghua University)</i> <i>Mingxing Xu (Tsinghua University)</i> <i>Thomas Fang Zheng (CSLT, Tsinghua University)</i></p>
P5-23-MLS (#251)	<p>Semi-Supervised Learning for Robust Speech Evaluation</p> <p><i>Huayun Zhang (ASTAR)</i> <i>Jeremy H. M. Wong (Institute for Infocomm Research)</i> <i>Geyu Lin (Agency of Science and Technology Research)</i> <i>Nancy Chen (Institute for Infocomm Research)</i></p>
P5-24-MLS (#264)	<p>GE2E-KWS: Generalized End-to-End Training and Evaluation for Zero-Shot Keyword Spotting</p> <p><i>Pai Zhu (Google)</i> <i>Jacob W Bartel (Google LLC)</i> <i>Dhruuv Agarwal (Google LLC)</i> <i>Kurt Partridge (Google)</i> <i>Hyun Jin Park (Google Inc.)</i> <i>Quan Wang (Google)</i></p>
P5-25-MLS (#312)	<p>A Simple HMM with Self-Supervised Representations for Phone Segmentation</p> <p><i>Gene-Ping Yang (The University of Edinburgh)</i> <i>Hao Tang (The University of Edinburgh)</i></p>
P5-26-MLS (#331)	<p>DASS: Distilled Audio State Space Models are Stronger and More Duration-Scalable Learners</p> <p><i>Saurabhchand Bhati (MIT)</i> <i>Yuan Gong (Massachusetts Institute of Technology)</i> <i>Leonid Karlinsky (MIT-IBM Watson AI Lab, IBM Research)</i> <i>Hilde Kuehne (University of Bonn)</i> <i>Rogério Feris (MIT-IBM Watson AI Lab, IBM Research)</i> <i>James Glass (Massachusetts Institute of Technology)</i></p>

P5-27-MLS (#337)	RAND: Robustness Aware Norm Decay for Quantized Neural Networks <i>David Qiu (Google)</i> <i>David Rim (Google)</i> <i>Shaojin Ding (Google)</i> <i>Oleg Rybakov (Google)</i> <i>Yanzhang He (Google)</i>
P5-28-MLS (#377)	SWIM: Short-Window CNN Integrated with Mamba for EEG-Based Auditory Spatial Attention Decoding <i>Ziyang Zhang (Tsinghua University)</i> <i>Andrew Thwaites (University College London)</i> <i>Alexandra Woolgar (University of Cambridge)</i> <i>Brian C.J. Moore (University of Cambridge)</i> <i>Chao Zhang (Tsinghua University)</i>

10:30-12:30 Challenge Session 5: FutureDial RAG (Venue: Lecture Hall)

12:30-14:00 Lunch

14:00-15:00 Invited Talk 3 (Venue: Lecture Hall)

Title: Towards Safe, Truly Open, and Factual Large Language Models

Speaker: *Prof Preslav Nakov, Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi*

Chair: *Dr Zhijian Ou*

15:00-15:30 Coffee Break

15:30-17:30 Panel Discussion

18:00-20:30 Gala Dinner

Day 4, Dec 5, Thursday

09:00-10:00 Keynote Speech 4 (Venue: Lecture Hall)

Title: A Theory of Unsupervised Speech Recognition

Speaker: *Prof. Mark Hasegawa-Johnson, University of Illinois at Urbana-Champaign*

Chair: *Prof Kong Aik Lee*

10:00-10:30 Coffee Break

10:30-12:30 Poster Session 6: Spoken Language Processing (Venue: Poster Area)

Chair: Dr Nan Yan

Poster ID (Paper ID)	Title and Authors
P3-28-SS04 (#411)	The Database and Benchmark for the Source Speaker Tracing Challenge 2024 <i>Ze Li (Wuhan University)</i> <i>Yuke Lin (Duke Kunshan University)</i> <i>Yao Tian (AI Center, OPPO)</i> <i>Hongbin Suo (AI Center, OPPO)</i> <i>Pengyuan Zhang (Institute of Acoustics, Chinese Academy of Sciences)</i> <i>Yanzhen Ren (Computer School of Wuhan University)</i> <i>Zexin Cai (Johns Hopkins University)</i> <i>Hirimitsu Nishizaki (University of Yamanashi)</i> <i>Ming Li (Duke Kunshan University)</i>
P6-01-SLP (#65)	Stutter-Solver: End-to-End Multi-Lingual Dysfluency Detection <i>Xuanru Zhou (Berkeley Speech Group)</i> <i>Cheol Jun Cho (UC Berkeley)</i> <i>Ayati Sharma (University of California, Berkeley)</i> <i>Brittany Morin (UCSF)</i> <i>David Baquirin (UCSF)</i> <i>Jet Vonk (UCSF)</i> <i>Zoe Ezzes (UCSF)</i> <i>Zachary Miller (UCSF)</i> <i>Boon Lead Tee (UCSF)</i> <i>Maria Luisa Gorno Tempini (UCSF)</i> <i>Jiachen Lian (University of California, Berkeley)</i> <i>Gopala Krishna Anumanchipalli (UC Berkeley)</i>
P6-02-SS09 (#417)	Domain Adaption and Unified Knowledge Base Motivate Better Retrieval Models in Dialog Systems with RAG <i>Huadong Lin (South China University of Technology)</i> <i>Yirong Chen (South China University of Technology)</i> <i>Wenyu Tao (South China University of Technology)</i> <i>Mingyu Chen (South China University of Technology)</i> <i>Xiangmin Xu (South China University of Technology)</i> <i>Xiaofen Xing (South China University of Technology)</i>
P6-03-SLP (#153)	SSAMBA: Self-Supervised Audio Representation Learning with Mamba State Space Model <i>Siavash Shams (Columbia University)</i>

	<p><i>Sukru Samet Dindar (Columbia University)</i> <i>Xilin Jiang (Columbia University)</i> <i>Nima Mesgarani (Columbia University)</i></p>
P6-04-SLP (#215)	<p>Speech-Copilot: Leveraging Large Language Models for Speech Processing via Task Decomposition, Modularization, and Program Generation</p> <p><i>Chun-Yi Kuan (National Taiwan University)</i> <i>Chih-Kai Yang (National Taiwan University)</i> <i>Wei-Ping Huang (National Taiwan University)</i> <i>Ke-Han Lu (National Taiwan University)</i> <i>Hung-Yi Lee (National Taiwan University)</i></p>
P6-05-SLP (#174)	<p>CTC-GMM: CTC-Guided Modality Matching for Fast and Accurate Streaming Speech Translation</p> <p><i>Rui Zhao (Microsoft)</i> <i>Jinyu Li (Microsoft)</i> <i>Ruchao Fan (Microsoft)</i> <i>Matt Post (Microsoft)</i></p>
P6-06-SLP (#223)	<p>Long-Form End-to-End Speech Translation via Latent Alignment Segmentation</p> <p><i>Peter Polák (Charles University)</i> <i>Ondrej Bojar (Charles University)</i></p>
P6-07-SLP (#272)	<p>Confidence Estimation for LLM-Based Dialogue State Tracking</p> <p><i>Yijyun Sun (University of Illinois, Urbana-Champaign)</i> <i>Suvodip Dey (University of Illinois, Urbana-Champaign)</i> <i>Dilek Hakkani-Tur (University of Illinois, Urbana-Champaign)</i> <i>Gokhan Tur (Amazon)</i></p>
P6-08-SLP (#278)	<p>The 2nd FutureDial Challenge: Dialog Systems with Retrieval Augmented Generation (FutureDial-RAG)</p> <p><i>Yucheng Cai (Tsinghua University)</i> <i>Si Chen (China Mobile Research)</i> <i>Yuxuan Wu (Tsinghua University)</i> <i>Yi Huang (China Mobile Research)</i> <i>Junlan Feng (China Mobile Research)</i> <i>Zhijian Ou (Tsinghua University)</i></p>
P6-09-SLP (#130)	<p>Zero-Shot Audio Topic Reranking Using Large Language Models</p> <p><i>Mengjie Qian (Cambridge University)</i> <i>Rao Ma (University of Cambridge)</i></p>

	<p><i>Aidan Liusie (University of Cambridge)</i> <i>Erfan Loweimi (University of Cambridge)</i> <i>Katherine M Knill (University of Cambridge)</i> <i>Mark Gales (University of Cambridge)</i></p>
P6-10-SLP (#59)	<p>Clean Label Attacks Against SLU Systems</p> <p><i>Henry Li Xinyuan (Johns Hopkins University)</i> <i>Thomas Thebaud (Johns Hopkins University)</i> <i>Sonal Joshi (Johns Hopkins University)</i> <i>Jesus Antonio Villalba (Johns Hopkins University)</i> <i>Najim Dehak (Johns Hopkins University)</i> <i>Sanjeev Khudanpur (Johns Hopkins University)</i></p>
P6-11-SLP (#168)	<p>WHISMA: A Speech-LLM to Perform Zero-Shot Spoken Language Understanding</p> <p><i>Mohan Li (Toshiba Europe Ltd.)</i> <i>Cong-Thanh Do (Toshiba Research Europe Ltd.)</i> <i>Simon Keizer (Toshiba Europe Ltd.)</i> <i>Youmna Farag (Toshiba Europe Ltd.)</i> <i>Svetlana Stoyanchev (Toshiba Europe Ltd.)</i> <i>Rama S Doddipatla (Toshiba Europe Ltd.)</i></p>
P6-12-SLP (#169)	<p>Improving Transducer-Based Spoken Language Understanding with Self-Conditioned CTC and Knowledge Transfer</p> <p><i>Vishal Sunder (The Ohio State University)</i> <i>Eric Fosler-Lussier (The Ohio State University)</i></p>
P6-13-SLP (#208)	<p>Self-Supervised Syllable Discovery Based on Speaker-Disentangled HuBERT</p> <p><i>Ryota Komatsu (Independent Researcher)</i> <i>Takahiro Shinozaki (Tokyo Institute of Technology)</i></p>
P6-14-SLP (#355)	<p>Just ASR + LLM? A Study on Speech Large Language Models' Ability to Identify and Understand Speakers in Spoken Dialogue</p> <p><i>Junkai Wu (University of Washington)</i> <i>Xulin Fan (University of Illinois at Urbana-Champaign)</i> <i>Bo-Ru Lu (University of Washington)</i> <i>Xilin Jiang (Columbia University)</i> <i>Nima Mesgarani (Columbia University)</i> <i>Mark A Hasegawa-Johnson (University of Illinois)</i> <i>Mari Ostendorf (University of Washington)</i></p>
P6-15-SLR (#5)	<p>Enhancing Open-Set Speaker Identification through Rapid Tuning with Speaker Reciprocal Points and Negative Samples</p>

	<p><i>Zhiyong Chen (Shanghai University)</i> <i>Zhiqi Ai (Shanghai University)</i> <i>Xinnuo Li (Shanghai University)</i> <i>Shugong Xu (Shanghai University)</i></p>
P6-16-SLR (#10)	<p>Spoofing-Aware Speaker Verification Robust Against Domain and Channel Mismatches</p> <p><i>Chang Zeng (National Institute of Informatics)</i> <i>Xiaoxiao Miao (Singapore Institute of Technology)</i> <i>Xin Wang (National Institute of Informatics)</i> <i>Erica Cooper (National Institute of Information and Communications Technology)</i> <i>Junichi Yamagishi (National Institute of Informatics)</i></p>
P6-17-SLR (#49)	<p>Adversarial Purification for Speaker Verification by Two-Stage Diffusion Models</p> <p><i>Yibo Bai (The University of Hong Kong)</i> <i>Zhang Xiaolei (Northwestern Polytechnical University)</i> <i>Xuelong Li (Institute of Artificial Intelligence (TeleAI), China Telecom Corp. Ltd.)</i></p>
P6-18-SLR (#103)	<p>Measuring Sound Symbolism in Audio-Visual Models</p> <p><i>Wei-Cheng Tseng (The University of Texas at Austin)</i> <i>Yi-Jen Shih (The University of Texas at Austin)</i> <i>David Harwath (The University of Texas at Austin)</i> <i>Raymond Mooney (The University of Texas at Austin)</i></p>
P6-19-SLR (#111)	<p>Meta-Learning Approaches for Improving Detection of Unseen Speech Deepfakes</p> <p><i>Ivan Kukanov (KLASS Engineering and Solutions)</i> <i>Janne Laakkonen (UEF)</i> <i>Tomi H. Kinnunen (University of Eastern Finland)</i> <i>Ville Hautamäki (University of Eastern Finland)</i></p>
P6-20-SLR (#135)	<p>On the Generation and Removal of Speaker Adversarial Perturbation for Voice-Privacy Protection</p> <p><i>Chenyang Guo (University of Science and Technology of China)</i> <i>Liping Chen (University of Science and Technology of China)</i> <i>Zhuhai Li (University of Science and Technology of China)</i> <i>Kong Aik Lee (The Hong Kong Polytechnic University)</i> <i>Zhen-Hua Ling (University of Science and Technology of China)</i> <i>Wu Guo (University of Science and Technology of China)</i></p>

<p>P6-21-SLR (#138)</p>	<p>Towards Quantifying and Reducing Language Mismatch Effects in Cross-Lingual Speech Anti-Spoofing</p> <p><i>Tianchi Liu (National University of Singapore)</i> <i>Ivan Kukanov (KLASS Engineering and Solutions)</i> <i>Zihan Pan (Institute for Infocomm Research (I2R), ASTAR, Singapore)</i> <i>Qiongqiong Wang (ASTAR)</i> <i>Hardik B Sailor (I2R, ASTAR, Singapore)</i> <i>Kong Aik Lee (The Hong Kong Polytechnic University)</i></p>
<p>P6-22-SLR (#141)</p>	<p>Enhancing Low-Resource Spoken Language Identification via Cross-Modality Retrieval and Cross-Lingual Text-to-Speech Synthesis</p> <p><i>Min Ma (Google DeepMind)</i> <i>Yuan Wang (Google)</i> <i>Kyle Kastner (Google)</i> <i>Isaac Caswell (Google)</i> <i>Charles Yoon (Google)</i> <i>Andrew Rosenberg (Google LLC)</i></p>
<p>P6-23-SLR (#148)</p>	<p>Recursive Attentive Pooling for Extracting Speaker Embeddings from Multi-Speaker Recordings</p> <p><i>Shota Horiguchi (NTT Corporation)</i> <i>Atsushi Ando (NTT Corporation)</i> <i>Takafumi Moriya (NTT Corporation)</i> <i>Takanori Ashihara (NTT Corp.)</i> <i>Hiroshi Sato (NTT)</i> <i>Naohiro Tawara (NTT)</i> <i>Marc Delcroix (NTT)</i></p>
<p>P6-24-SLR (#259)</p>	<p>PDAF: A Phonetic Debiasing Attention Framework for Speaker Verification</p> <p><i>Massa Baali (CMU)</i> <i>Abdulhamid Aldoobi (Carnegie Mellon University)</i> <i>Hira Dharmyal (Carnegie Mellon University)</i> <i>Rita Singh (Carnegie Mellon University)</i> <i>Bhiksha Raj (Carnegie Mellon University)</i></p>
<p>P6-25-SLR (#356)</p>	<p>INX-SpeakerHub: A 2000-Hour Indian Multilingual Speaker Identification Corpus</p> <p><i>Metilda Sagaya Mary NJ (Indian Institute of Technology Madras)</i> <i>S Umesh (IIT Chennai)</i></p>
<p>P6-26-DIA (#100)</p>	<p>Resource-Efficient Adaptation of Speech Foundation Models for Multi-Speaker ASR</p>

	<p><i>Weiqing Wang (NVIDIA)</i> <i>Kunal Dhawan (NVIDIA)</i> <i>Tae Jin Park (NVIDIA)</i> <i>Krishna C Puvvada (NVIDIA)</i> <i>Ivan Medennikov (NVIDIA)</i> <i>Somshubra Majumdar (NVIDIA)</i> <i>He Huang (NVIDIA)</i> <i>Jagadeesh Balam (NVIDIA)</i> <i>Boris Ginsburg (NVIDIA)</i></p>
P6-27-SS01 (#405)	<p>Exploring Self-Supervised Representations for Text-Dependent Speaker Verification</p> <p><i>Sankala Sreekanth (Indian Institute of Technology Hyderabad (IITH))</i></p>
P6-28-SS04 (#147)	<p>Distillation-Based Feature Extraction Algorithm for Source Speaker Verification</p> <p><i>Xinlei Ma (Tianjin University)</i> <i>Wenhuan Lu (Tianjin University)</i> <i>Ruiteng Zhang (Tianjin University)</i> <i>Junhai Xu (Tianjin Key Laboratory of Cognitive Computing and Application, College of Intelligence and Computing, Tianjin University)</i> <i>Xugang Lu (NICT)</i> <i>Jianguo Wei (School of Computer Software, Tianjin University, Tianjin, China)</i></p>
P6-29-SS04 (#224)	<p>Speaker Contrastive Learning for Source Speaker Tracing</p> <p><i>Qing Wang (Northwestern Polytechnical University)</i> <i>Hongmei Guo (Northwestern Polytechnical University)</i> <i>Jian Kang (Institute of Artificial Intelligence (TeleAI), China Telecom)</i> <i>Mengjie Du (China Telecom)</i> <i>Jie Li (Institute of Artificial Intelligence (TeleAI), China Telecom)</i> <i>Zhang Xiaolei (Northwestern Polytechnical University)</i> <i>Lei Xie (NWPU)</i></p>

10:30-12:30 Challenge Session 6: GenASR challenge (Venue: Lecture Hall)

12:30-13:00 Closing Ceremony (Venue: Lecture Hall)