## The VoiceMOS Challenge 2023: Zero-shot Subjective Speech Quality Prediction for Multiple Domains

Erica Cooper, Wen-Chin Huang, Yu Tsao, Hsin-Min Wang, Tomoki Toda, Junichi Yamagishi Nagoya University; National Institute of Informatics, Japan; Academia Sinica, Taiwan

## The VoiceMOS Challenge (VMC) VoiceMOS Challenge 2022 Focuses on automatic quality MM assessment for synthetic speech \*\*\* Attracted more than 20 participants MOS Challe Main track: the BVCC dataset Mean opinion score (MOS) test: Subjective Rate quality of individual samples. Large-scale re-evaluation of TTS & VC samples from 2008 test -MA Best system: .979 system-level SRCC -~\\M^ ∮→ 🗙 🗙 🗙 $\rightarrow$ Current technology performs well in an in-domain setting Out-of-domain track: Blizzard Challenge 2019 Chinese TTS: Small amount of labeled data (136 samples) Data-driven MOS prediction Best system: .975 system-level SRCC (mostly based on deep learning) $\rightarrow$ Current technology performs well in a fine-tuning setting Obiective -MN~ -{ 🔲 → 🗙 ★ ★ assessment Still far from a truly generalizable, zero-shot predictor! Generalizable? Track 1: French text-to-speech Track 2: Singing voice conversion SINGING VOICE CONVERSION É In collaboration with the Voice Conversion Challenge In collaboration with Blizzard Challenge (a VC challenge founded in 2016) (TTS challenge founded in 2008); Theme this year: French Theme this year: singing voice conversion Two sub tracks: (1a) Speaker-dependent; (1b) Speaker-adaptive Datasets Track 3: noisy/enhanced speech No official training set for tracks 1 & 2! Dataset: TMHINT-QI (Mandarin) Samples # ratings Track Туре Lang Systems Training set per system per sample 4 noise types (babble, street, pink, and white) at 4 SNR levels (-2, 0, 2, and 5) Track 1a Hub: 21 42 TTS 15 Fre 5 enhancement (SE) models: KLT, MMSE, FCN, DDAE, and Transformer Track 1b Spoke: 17 34 In-dom: 25 Singing Testing set Track 2 Eng 80 6 VC Cross-dom: 24 Same noise generation process as in training set ٠ Noisy & 5 SE models: MMSE, FCN, Trans, DEMUCS, and CMGAN Track 3 Chi 97 20 5.3 enhanced Results VoiceMOS Challenge 2023 Results 1.0 Track 1a • Raw waveform input Spectrogram input Latent feature inpu Track 1b Track 2 Track 3 0.8 0.8 ZZ Post-e SRCC KTAU 9.0 0.4 S Svs. 4 4 0.2 0.2 0.0 5-P17-SS07 5-V9-SS07 5-V8-SS07 0.0 т02 ●▲◆ тоз то́4 ●▲ т05 ▲■ т09 ▲■ во1 • во́2 то1 то́6 ●∎ то́8 тіо тiı The goal of having one model that can predict MOS well for different **Best vs. Average scores** domains has not yet been reached. Best score for each track Average score for each track **Domain mismatch** between singing voice conversion and text-to-speech synthesis was **not as large** as we expected. Track SRCC KTAU Team Avg SRCC Avg KTAU The **difference** in predictability between **speaker-dependent** and 0.790.60 T06 0.57 0.42 1a speaker-adaptive French TTS was surprising. 0.91 T05 0.50 0.39 0.75 1b 2 0.87 0.69 T03 0.67 0.50 Approaches using **listener information** and a **mix of training datasets** 3 0.95 0.83 T02 0.63 0.49 tended to be more successful. Track 3 has highest best score Track 2 has highest average score We can **observe a gap** between the case where a **small amount of in**domain training data is available, and the completely zero-shot setting. Feedback Future directions Challenge HP Keep on improving zero-shot ability Unfamiliarity with French TTS, Singing voice and speech enhancement is a difficulty Moving from MOS to other listening tests (MUSHRA, AB preference tests, ...) High listener variance in Track 3 More fine-grained instructions (expressivity, etc.) Not sure what is a proper training set for tracks 1 & 2

