

# The VoiceMOS Challenge 2024: Beyond Speech Quality Prediction

P4-26-SS05  
(#396)

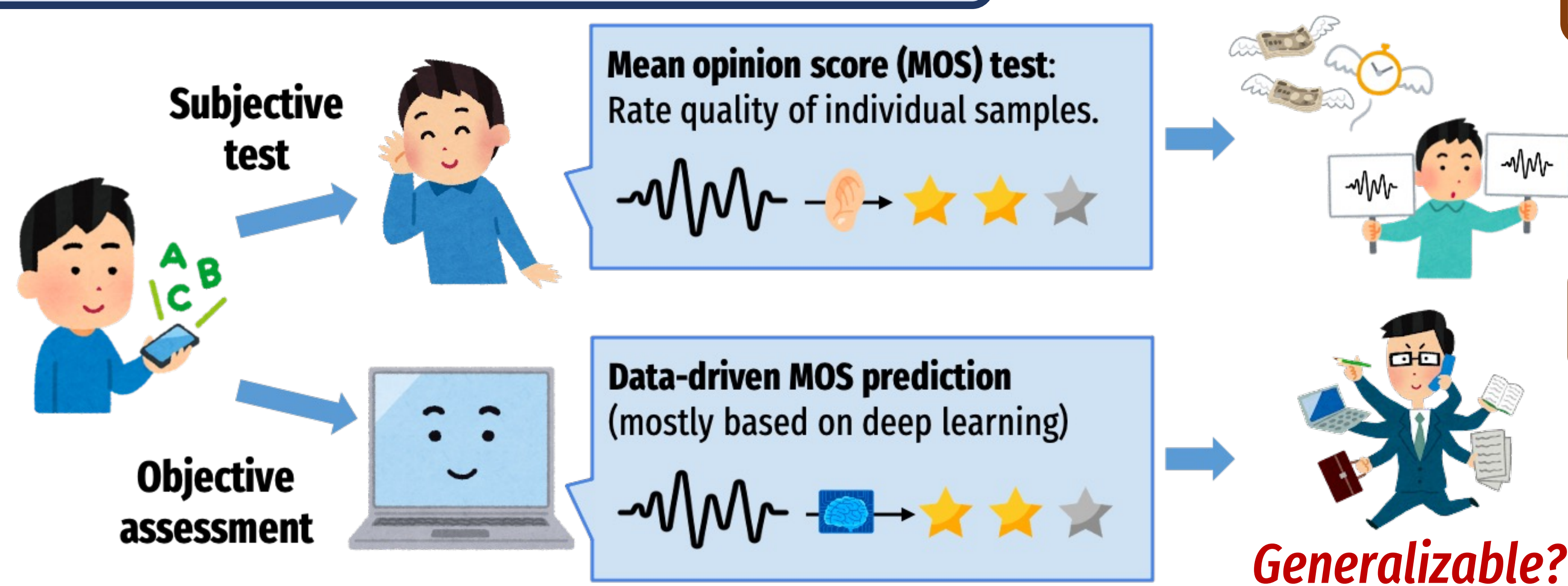
Wen-Chin Huang<sup>1</sup>, Szu-Wei Fu<sup>2</sup>, Erica Cooper<sup>3</sup>, Ryandhimas E. Zezario<sup>4</sup>,

Tomoki Toda<sup>1</sup>, Hsin-Min Wang<sup>4</sup>, Junichi Yamagishi<sup>5</sup>, Yu Tsao<sup>4</sup>

<sup>1</sup>Nagoya University, JP <sup>2</sup>NVIDIA Taiwan <sup>3</sup>NICT, JP <sup>4</sup>Academia Sinica, Taiwan <sup>5</sup>NII, JP



## The VoiceMOS Challenge (VMC) series



2022

- Focus: synthetic speech, supervised setting
- Datasets: BVCC dataset & Blizzard Challenge (BC) '19
- Large-scale re-evaluation of TTS & VC samples since '08
- Best system: .979/.975 system-level SRCC  
→ Performs well in the supervised setting

2023

- Focus: zero-shot setting
- Tracks: Blizzard challenge (French TTS), Singing Voice Conversion Challenge, clean/noisy/enhanced speech
- Result 1: gap between supervised & zero-shot setting
- Result 2: no consistent performance across all tracks

## Track 1: MOS prediction for "zoomed-in" systems

- Motivation: evaluate synthetic systems of high-quality
- New listening tests using the top 50%, 25%, 12% systems in BVCC
- 50% → validation set; 25% & 12% → test set

## Track 2: MOS prediction for singing voice

- A newly collected dataset named SingMOS: natural singing voices, vocoder analysis-synthesis, singing voice synthesis/conversion samples
- Mandarin & Japanese, 16kHz, 35 systems, 2000/544/645 samples

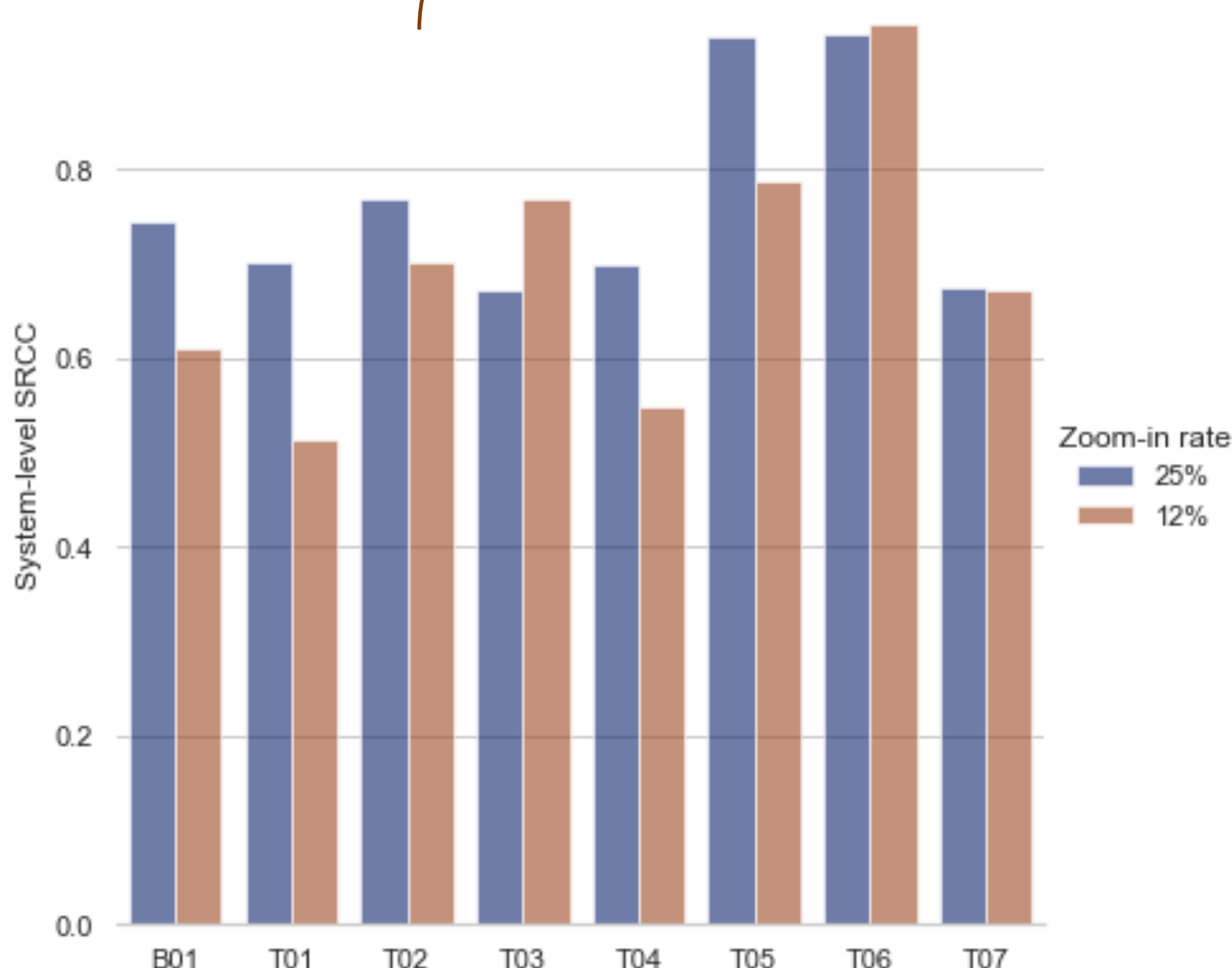
## Track 3: semi-supervised MOS prediction for clean/noisy/enhanced speech

- Setting: very limited amount of training data & zero-shot setting
- Train/valid set: UDASE task of 7<sup>th</sup> CHiME, 60/40 samples (real noisy samples)  
Test set: VoiceBank-DEMAND, 4 noise types, 5 enhancement systems, 280 samples (artificial samples)
- Beyond quality: speech signal quality (SIG), background intrusiveness (BAK), overall quality (OVRL)

## Participants, baselines

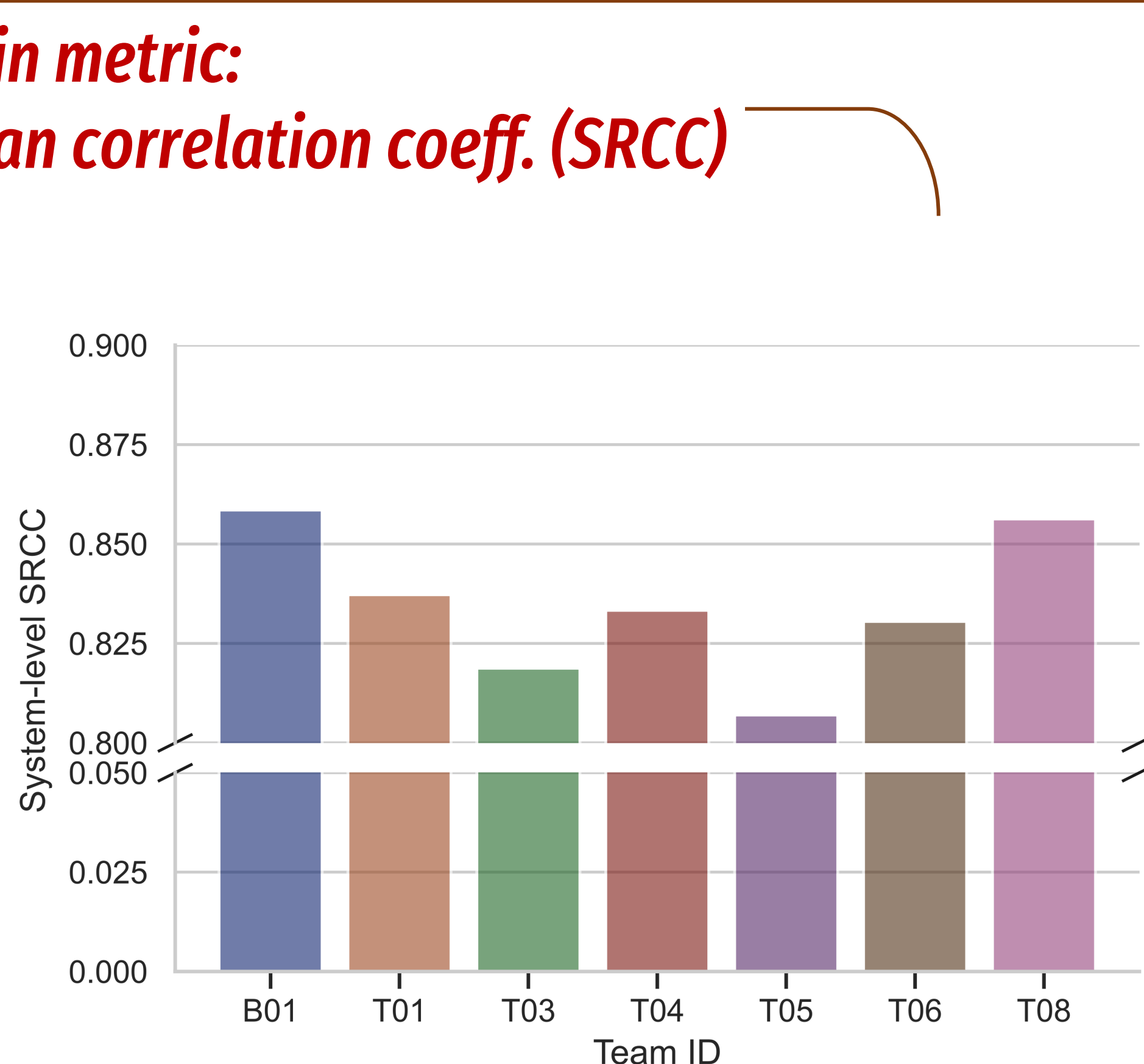
- 8 teams (5 academia, 3 industry)
- Baselines for tracks 1 & 2: SSL-MOS
- Baseline for track 3: VQScore

Main metric:  
system-level Spearman correlation coeff. (SRCC)



### Track 1:

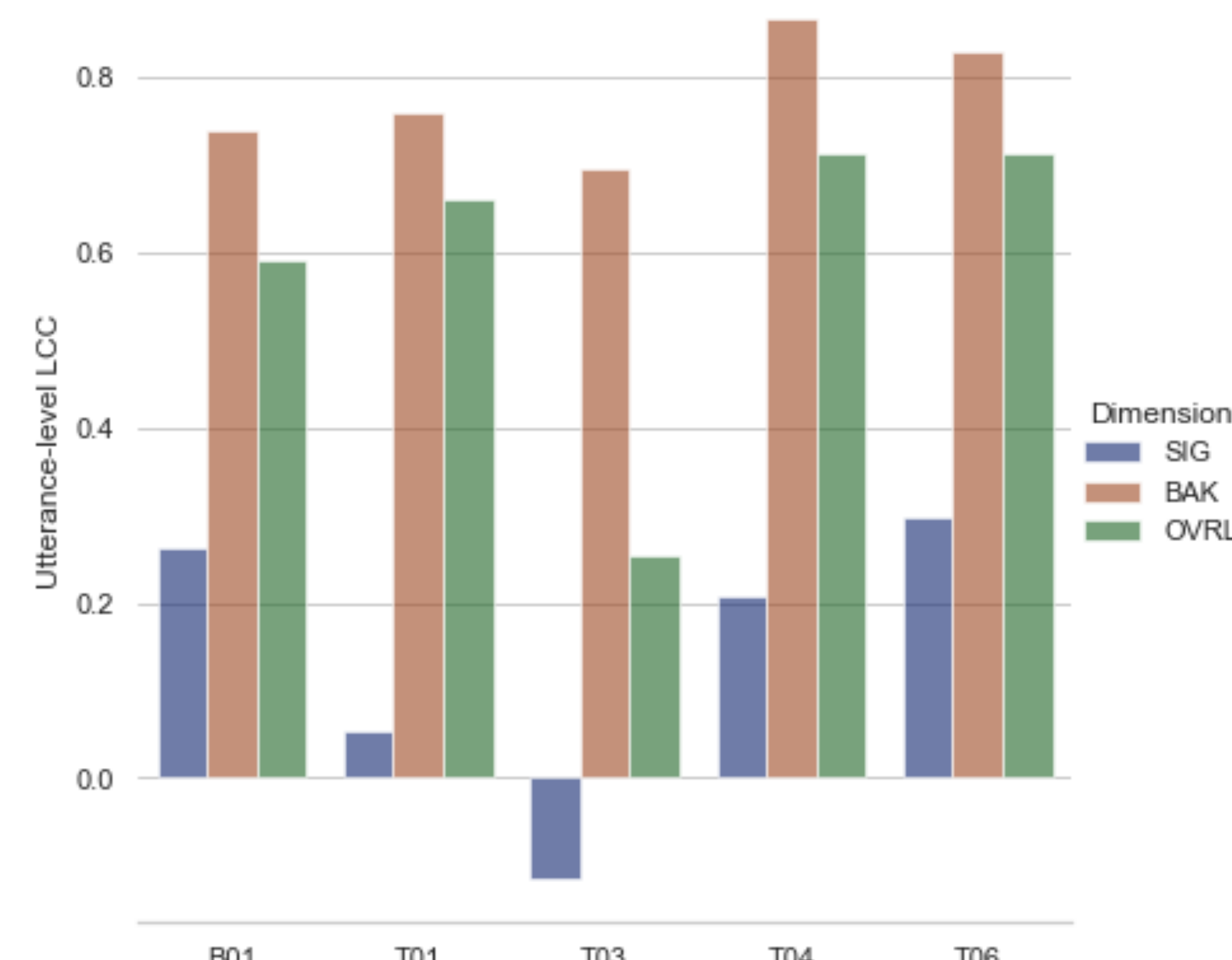
- 12% is harder than 25%
- Baseline (B01) ranked 4<sup>th</sup>/6<sup>th</sup> in 25%/12%  
→ participants have advanced
- Top systems: T05 & T06



### Track 2:

- No team outperformed the baseline (B01)
- Differences were small
- T06 ranked 1<sup>st</sup> in all utterance-level metrics

Main metric:  
utterance-level linear correlation coeff. (LCC)



### Track 3:

- Baseline (B01) was outperformed
- SIG is the most difficult to predict
- No team excelled all aspects
- T06: 1<sup>st</sup> in SIG & OVRL; T04: 1<sup>st</sup> in BAK

## Top system: T06

- Performed remarkably well in all three tracks.
- Improved version of RAMP: equipping a parametric model (e.g., SSL-MOS) with a non-parametric head based on kNNs.
- Was shown to generalize well to unseen data.

## Top system: T05 (P4-28-SS05 (#407))

- Top system in track 1.
- SSL feature + mel spectrogram (EfficientNetV2 encoder).
- Conducted own listening test.

## Top system: T08 (P4-27-SS05 (#406))

- Top system in track 2.
- SSL feature + pitch histogram.

## Top system: T04

- Top system in track 3.
- Trained separate models for BAK and SIG prediction.  $OVRL = (BAK + SIG) / 2$ .
- BAK predictor: pre-trained to predict SNR of simulated noisy speech samples.
- SIG predictor: pre-trained to predict spoofed and natural samples from ASVspoof 2019.
- Both are fine-tuned on the provided training data.

## Future directions

- Modern-day speech synthesis systems
- More diverse speech types
- Beyond speech: music, environmental sounds

## Challenge HP

