04

Neural beamforming for tracking moving sources

Keeping listening while talker is moving

Abstract

Speech enhancement technology is crucial for machines to correctly recognize human speech in noisy environments because it allows extracting only the voices we want to listen to from the background noise. This research introduces a novel beamforming technique that tracks the speaker's movement and keeps extracting the target speaker's voice, even when the speaker is moving while talking. Beamforming requires information on the direction of arrival of the target source and noise signals (spatial information). In this research, we consider the problem of estimating time-varying spatial information as a problem of segmenting speech according to the speaker's movement and propose a novel framework for solving this problem using deep learning. This framework enables accurate estimation of spatial information even when the target speaker is moving. With this framework, we aim for a future in which people and machines can interact more naturally in any situation, such as when multiple speakers talk freely while walking around.



References

[1] T. Ochiai, M. Delcroix, T. Nakatani, S. Araki, "Mask-Based Neural Beamforming for Moving Speakers With Self-Attention-Based Tracking," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, Vol. 31, pp. 835-848, 2023.

Contact

Tsubasa Ochiai, Signal Processing Research Group, Media Information Laboratory