# 08 Cleaning-up speech from noisy, reverberant recordings
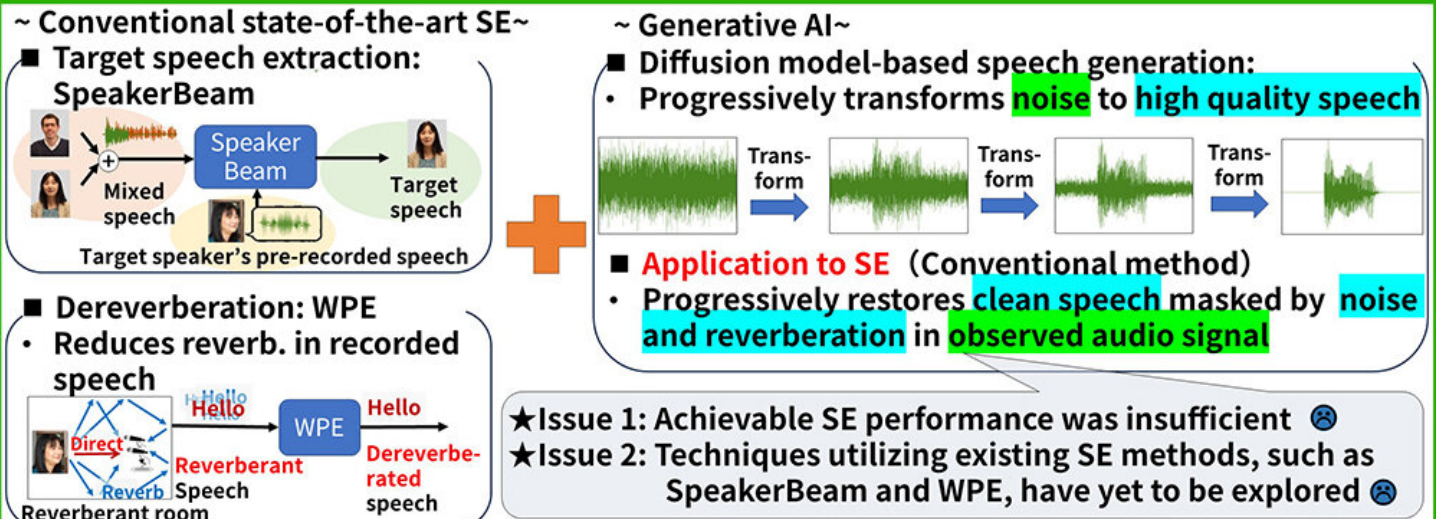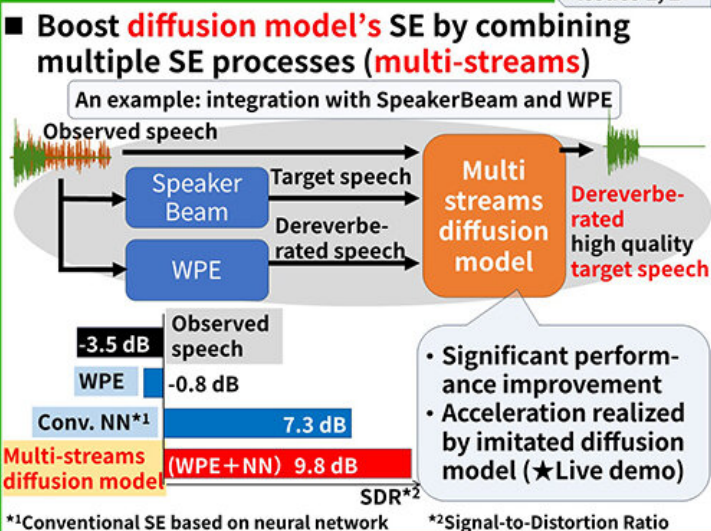
## Abstract

Recent advancements in generative AI have enabled high-quality speech enhancement. This research introduces a speech enhancement method that utilizes the diffusion model, one of the most powerful generative AI models, to effectively remove noise and reverberation from speech recordings. Our approach integrates multiple conventional speech enhancement techniques into a diffusion model-based framework, significantly improving performance. Additionally, we are the first in the world to demonstrate that averaging multiple outputs from the diffusion model, a technique we refer to as "ensemble inference", greatly enhancing performance. In the future, this technology will enable high-quality speech recording even in noisy environments, making voices sound as if recorded in a studio. This advancement is expected to greatly enhance various speech applications, such as collecting high-quality audio data in everyday environments and enabling more comfortable remote meetings.

## References

[1] N. Kamo, M. Delcroix, T. Nakatani, "Target speech extraction with conditional diffusion model," in *Proc. INTERSPEECH,* pp. 176-180, 2023.
[2] T. Nakatani, N. Kamo, M. Delcroix, S. Araki, "Multi-stream diffusion model for probabilistic integration of model-based and data-driven speech enhancement," in *Proc. IWAENC,* pp. 65-69, 2024.

## Contact

Naoyuki Kamo, Signal Processing Research Group, Media Information Laboratory