

# NAVER CLOVA SUBMISSION TO THE THIRD DIHARD CHALLENGE

*Hee-Soo Heo<sup>1</sup>, Jee-weon Jung<sup>1,2</sup>, Youngki Kwon<sup>1</sup>, You Jin Kim<sup>1</sup>,  
Jaesung Huh<sup>3</sup>, Joon Son Chung<sup>1</sup>, Bong-Jin Lee<sup>1</sup>*

<sup>1</sup>Naver Corporation, South Korea

<sup>2</sup>School of Computer Science, University of Seoul, South Korea

<sup>3</sup>Visual Geometry Group, Department of Engineering Science, University of Oxford, UK

## ABSTRACT

This report describes Naver Clova’s submission to the third DIHARD speech diarization challenge. Our system consists of following subsystems: speech activity detection, overlapped speech detection, speaker embedding extraction, feature enhancement, and clustering. Main improvements of our submitted system over existing diarization systems are feature enhancement based on the dimensionality reduction and overlap detection system. We reduce the dimensionality using an auto-encoder for each utterance to adapt speaker representations for the clustering system, and then perform attention-based aggregation, which we developed to enhance the clustering result. We also use some variants of CRNN based overlapped speech detection networks and their ensemble to further reduce the missed detection of overlapped speech regions. The submitted system achieves competitive performance in the DIHARDIII\_Task1\_CORE, which shows 14.96%, and 15.40% diarization error rate for the development and the evaluation sets, respectively.

***Index Terms***— Speaker diarisation.

## **1. REFERENCES**