

The Academia Sinica Systems of Speech Recognition and Speaker Diarization for the CHiME-6 Challenge

Hung-Shin Lee¹, Yu-Huai Peng¹, Pin-Tuan Huang¹, Ying-Chun Tseng², Chia-Hua Wu¹, Yu Tsao², Hsin-Min Wang¹

¹Institute of Information Science, Academia Sinica, Taiwan

²Research Center for Information Technology Innovation, Academia Sinica, Taiwan

hungshinlee@gmail.com

Abstract

This paper describes the Academia Sinica systems for the tracks of multiple-array ASR (Track 1) and diarization+ASR (Track2) in the 6th CHiME Challenge. For Track 1, we take a different approach from the official baseline to preprocess the Kinect data and derive the state-level alignment. In addition, we develop two LF-MMI-based acoustic models, the discriminative autoencoders (DcAE) and the feature-enhanced acoustic model (FEAM), which consider feature-level regularization and enhancement, respectively. For Track 2, we propose a new CNN-based training scheme, which develops speech representations by expanding the data into a set of segments, each of which contains more than one speaker. In training, a soft label is applied to each segment based on the speaker occupation ratio, and the standard cross entropy loss is used. In the evaluation set, our best system for Track 1 (Category A) achieves 46.8% WER, slightly better than the baseline performance (51.4%). For Track 2 (Category A), our system is also superior to the baseline while using the same TDNN-based acoustic model. The DER, JER, and WER are relatively improved by 13.24%, 12.60%, and 6.57%, respectively.

1. System Descriptions

We describe our systems for both tracks in the 6th CHiME Challenge (CHiME-6). For details of the CHiME-6 datasets and tasks, please refer to the official website¹ and [1].

1.1. Track 1: ASR

The training process of our ASR system is divided into two parts, front-end data processing and back-end acoustic modeling. As shown in the upper part of Figure 1, we first used the worn set and the Kinect set to train the GMMs. The worn set comes from the L and R channels in the worn microphone data, and is combined with the simulated reverberant speech using RIRs and point-source noises [2]. In the baseline program², the Kinect set consists of 400k utterances randomly selected from all Kinect channels without any enhancement. Our Kinect set, instead, comes from 1) all the first channel utterances of the Kinect data and 2) the corresponding enhanced utterances, where all channels with time annotations were passed to the front-end of weighted prediction error (WPE), guided source separation (GSS), and BeamformIt (BF) [3, 4, 5].

Following the model structure and training steps of the baseline program², we first created the phone alignment for the

¹<https://chimechallenge.github.io/chime6/overview.html>

²https://github.com/kaldi-asr/kaldi/tree/master/egs/chime6/s5_track1/

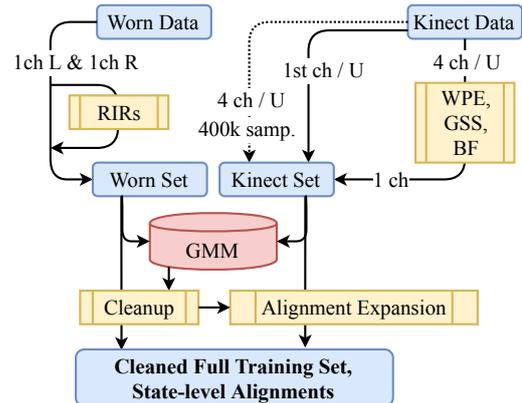


Figure 1: Flowchart of data preparation in our system, where “U” and “ch” denote “Kinect” and “channel”, respectively. The dotted line path is the method of forming the Kinect set in the baseline program.

worn set based on the GMMs, and performed a data cleanup procedure. We then created the alignment and lattice for the complete training set for the NN-based acoustic models (AMs) by copying the alignment of the corresponding L channel in the worn set, i.e., the alignment expansion in [6].

To train the NN-based AMs, the training set was augmented by two data augmentation techniques, namely speed perturbation and volume perturbation. Bandpass perturbation [6] was not successful in our experiments.

The architectures of our two newly proposed AMs are depicted in Figure 2. The first AM is discriminative autoencoders (DcAEs) [7], which attempts to effectively separate the phonetic part (P-Code) and the residual part (R-Code) in the embedding space. In this challenge, we not only corrected several minor mistakes in our previous implementation, but also upgraded its structure from “nnet3” to “chain”. In this way, the LF-MMI criterion, the cross-entropy loss, and the mean squared error can be optimized simultaneously by Kaldi’s training procedures.

The second AM is the feature-enhanced acoustic model (FEAM) as shown in Figure 2 (b). In FEAM-U, “-U” means that the U-Net is used. There are also two kinds of output layers, one is the phone-state scores for the LF-MMI criterion and the cross-entropy loss, and the other is the generated acoustic features. The acoustic features generated by the feature-enhanced networks (FENs) are expected to be close to the corresponding worn features during training. That is, we assume the worn set is almost clean, so that FENs can play a role in further enhancing the Kinect features.

In summary, we used five AMs, including DcAE-B, DcAE-U, FEAM-U, TDNN-F, and RBiLSTM [8, 6, 2]. All of them

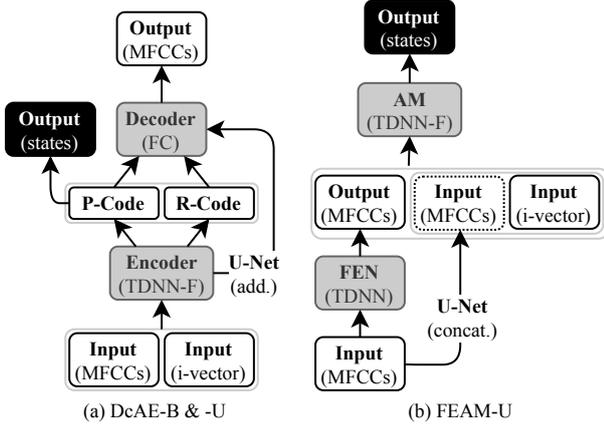


Figure 2: Our NN-based acoustic models. In (a), FC denotes the fully-connected layers, and U-Net is used in DcAE-U but not in DcAE-B. In (b), FEN denotes the feature-enhanced networks using the TDNN layers.

were trained on the “chain” structure using the Kaldi Toolkit, with the input combining 40-dimensional MFCCs and the 100-dimensional i-vector.

In the decoding phase, all Kinect channels were processed by WPE, GSS, and BF to form a single-channel utterance. Finally, we used the N-best ROVER method to combine the results from different AMs [9].

1.2. Tack 2: Diarization + ASR

For Track 2, we basically followed the baseline program³, including the constitution of training set, dereverberation procedures, speech activity detection (SAD), and the back-end of PLDA and AHC. In the baseline program, BF is used to combine and enhance all channels of each Kinect into one channel. However, in our system, some failures for unknown reasons occurred when BF was performed on all channels of all Kinects. Therefore, we tried all possible combinations of channels and selected the set that contains the most compatible channels.

A typical speaker diarization system is composed of two components, a speaker model and a back-end processor, working for extraction and clustering of speaker representations, respectively. The main weakness of most speaker models might be the incompetence to discriminate short-duration segments, e.g., less than 2 seconds, and the ineffectiveness to extract a reliable speaker embedding when a segment contains more than one speaker. Speaker representation is crucial to speaker diarization especially when segment clustering is performed. Therefore, we propose a new training scheme to develop speaker representations by randomly augmenting the training data with segments that contain more than one speaker. That is, we attempt to produce one or more “speaker change” in each mini-batch while training. Thereinto, a soft label was applied to each segment (sample) based on the speaker occupation ratio and the standard cross entropy loss was used. Take Figure 3 for example, if the ground truths of samples 1 and 2 in the traditional case are [1, 0, 0] and [0, 1, 0], they can be [3/8, 5/8, 0] and [2/8, 3/8, 3/8], respectively, in our proposed case. The ratio of the number of multi-speaker segments to the number of single-speaker segments is about 13.64%.

To build a speaker model, we employed the CNN-based

³https://github.com/kaldi-asr/kaldi/tree/master/egs/chime6/s5_track2/

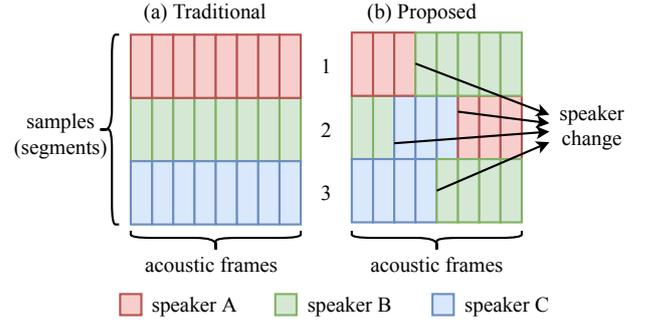


Figure 3: Illustration of the traditional training scheme and our proposed training scheme with respect to the speaker distribution in a mini-batch. Each sample (segment) contains 8 acoustic frames and the mini-batch size is 3.

ResNet-34 architecture, where the feature kind, training hyper-parameters, specification of layers, aggregation type, and loss function are almost the same as the baseline system described in [10]. The only difference is that we added one additional res-layer with 256 channels and 3 blocks to the original model in order to extract the 256-dimensional speaker embeddings. The VoxCeleb-2 Corpus was used for training the speaker model. The best model checkpoint was determined by JERs of the CHiME-6 development set.

The initial speaker label for each segment was given through clustering the speaker embeddings. Resegmentation was subsequently performed with variational Bayes (VB) diarization [11], where a 2048-component UBM-GMM with diagonal covariance matrices and 400 eigenvoice bases were trained in advance with 30-dimensional MFCCs. Moreover, the initial speaker label was used for initialization in the VB diarization model. The tunable parameters, such as the minimum duration, loop probability, downsampling factor, and maximum number of iteration, were determined by the development set, and were set to 1, 0.998, 1, and 1, respectively.

2. Experiment Results

The results for the development and evaluation sets are presented in Tables 1 and 2.

Table 1: WERs (%) for Track 1 and Track 2 (Category A only).

Model	Track 1		Track 2	
	Dev	Eval	Dev	Eval
Baseline	51.32	51.36	84.25	77.94
TDNN-F	50.12	49.36	75.89	73.68
RBiLSTM	52.43	50.26	76.90	73.39
DcAE-B	50.12	49.68	75.90	73.66
DcAE-U	49.86	49.63	75.78	73.54
FEAM-U	53.47	52.70	78.70	76.20
ROVER	47.28	46.82	74.36	71.56

Table 2: Results for Track 2. The acoustic models are the same.

Model	Dev			Eval		
	DER	JER	WER	DER	JER	WER
Baseline	63.42	70.83	84.25	68.20	72.54	77.94
Proposed	56.77	60.62	75.57	59.17	63.40	72.82

3. References

- [1] J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," in *Proc. Interspeech*, 2018.
- [2] V. Manohar, S.-j. Chen, Z. Wang, Y. Fujita, S. Watanabe, and S. Khudanpur, "Acoustic modeling for overlapping speech recognition: JHU CHiME-5 challenge system," in *Proc. ICASSP*, 2019.
- [3] L. Drude, J. Heymann, C. Boeddeker, and R. Haeb-Umbach, "NARA-WPE: A Python package for weighted prediction error dereverberation in Numpy and Tensorflow for online and offline processing," *Proc. 13th ITG-Symposium*, 2018.
- [4] C. Boeddeker, J. Heitkaemper, J. Schmalenstroerer, L. Drude, J. Heymann, and R. Haeb-umbach, "Front-end processing for the CHiME-5 dinner party scenario," in *Proc. Workshop on Speech Processing in Everyday Environments (CHiME)*, 2018.
- [5] M. Kitzka, W. Michel, C. Boeddeker, J. Heitkaemper, T. Menne, R. Schlüter, H. Ney, J. Schmalenstroerer, L. Drude, J. Heymann, and R. Haeb-Umbach, "The RWTH/UPB system combination for the CHiME 2018 workshop," in *Proc. Workshop on Speech Processing in Everyday Environments (CHiME)*, 2018.
- [6] N. Kanda, R. Ikeshita, S. Horiguchi, Y. Fujita, K. Nagamatsu, X. Wang, V. Manohar, N. E. Yalta Soplín, M. Maciejewski, S.-J. Chen, A. S. Subramanian, R. Li, Z. Wang, J. Naradowsky, L. P. Garcia-Perera, and G. Sell, "The Hitachi/JHU CHiME-5 system: Advances in speech recognition for everyday home environments using multiple microphone arrays," in *Proc. Workshop on Speech Processing in Everyday Environments (CHiME)*, 2018.
- [7] P.-T. Huang, H.-S. Lee, S.-S. Wang, K.-Y. Chen, Y. Tsao, and H.-M. Wang, "Exploring the encoder layers of discriminative autoencoders for LVCSR," in *Proc. Interspeech*, 2019.
- [8] D. Povey, G. Cheng, Y. Wang, K. Li, H. Xu, M. Yarmohamadi, and S. Khudanpur, "Semi-orthogonal low-rank matrix factorization for deep neural networks," in *Proc. Interspeech*, 2018.
- [9] J. G. Fiscus, "A post-processing system to yield reduced word error rates: Recognizer Output Voting Error Reduction (ROVER)," in *Proc. ASRU*, 1997.
- [10] W. Cai, J. Chen, and M. Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system," in *Proc. Odyssey*, 2018. [Online]. Available: http://www.isca-speech.org/archive/Odyssey_2018/abstracts/26.html
- [11] M. Diez, L. Burget, and P. Matejka, "Speaker diarization based on Bayesian HMM with eigenvoice priors," in *Proc. Odyssey*, 2018.