

A Novel End-to-End CAPT System for L2 Children Learners

Anonymous ACL submission

Abstract

001 Recently, Conformer-based model shows
002 promising results in automatic speech recog-
003 nition (ASR) task. There still is a dearth
004 of research on Conformer based model
005 for computer-assisted pronunciation learning
006 (CAPT) system. In this paper, a Conformer-
007 based CAPT system is introduced to provide
008 the mispronunciation detection and diagno-
009 sis. We apply the Conformer as the main pro-
010 nunciation error detection model in phoneme
011 level since superior phoneme recognition per-
012 formance. Then, the features, including the Log
013 Phone Posterior (LPP), the Log Posterior Ratio
014 (LPR) and some other features, extracted from
015 the Conformer decoder, are trained by a XG-
016 Boost model to predict phoneme and sentence
017 level scores labeled by experts. Both results on
018 open datasets and our internal Chinese children
019 data demonstrate that the Conformer-based sys-
020 tem, which has smaller model size and detailed
021 diagnosis, achieves better performance com-
022 pared with neural network (NN)-based system.

023 1 Introduction

024 The computer-aided pronunciation training (CAPT)
025 application plays an important role in the computer-
026 aided language learning (CALL) task especially for
027 the second language (L2) learners. The mispronun-
028 ciation detection and diagnosis (MD&D) module
029 is the core of CAPT system, providing individ-
030 ualized pronunciation evaluation and associated
031 phone-level diagnosis.

032 A classic method for MD&D module is pronun-
033 ciation scoring using the goodness of pronuncia-
034 tion (GOP) (Witt, 1999; Witt and Young, 2000)
035 which is based on the confidence measures of the
036 acoustic model (AM). GOP employs the ratio be-
037 tween the likelihood of canonical and the most
038 likely pronounced phones computed by the align-
039 ments on given text and the voice of learner. The
040 hybrid deep neural network-hidden markov models
041 (DNN-HMM) architecture is the mostly used in the

MD&D system (Hu et al., 2015; Zhang et al., 2021)
with more accurate measures using some discrimi-
native training algorithms. Some GOP variations
(Sudhakara et al., 2019; Wana et al., 2020) are also
introduced in the recent years showing good corre-
lates with human/expert labels.

042 However, the GOP based methods can only pro-
043 vide the pronunciation scores without insertion er-
044 rors in the pronunciation. It is hard to detect in-
045 serted phonemes/words and give more diagnostic
046 details for the language learners. To deal with these
047 shortcomings, some end-to-end (E2E) based meth-
048 ods (Feng et al., 2020; Yan et al., 2020; Fu et al.,
049 2021; Jiang et al., 2021) are introduced employing
050 a free-phone recognition system with connectionist
051 temporal classification (CTC) (Paterlini-Brechet
052 and Benali, 2007; Leung et al., 2019) or hybrid
053 CTC- Attention model (Watanabe et al., 2017; Yan
054 et al., 2020). With the development of automatic
055 speech recognition (ASR), Conformer (Gulati et al.,
056 2020) based E2E models achieves state-of-the-art
057 (SOTA) performance by combining the convolu-
058 tion neural networks (CNN) and Transformers, to
059 model both local and global dependencies of an
060 audio sequence in a parameter-efficient way. Nev-
061 ertheless, there is still a lack of CAPT system de-
062 signed by Conformer model to deal with both MD
063 and diagnosis. So it is necessary to compute a score
064 based on the recognized results to meet the require-
065 ments of the exercises and the examinations in the
066 language education.

067 In this paper, we propose a Conformer-based
068 CAPT system focused on MD&D. Firstly, a Con-
069 former model is trained in the phonetic level. Then,
070 with the force-alignment and the posterior probabili-
071 ty given by the trained model, features such as the
072 Log Phone Posterior (LPP), the Log Posterior Ratio
073 (LPR) (Hu et al., 2015) and some other features are
074 extracted. Finally, these features are trained with a
075 XGBoost model to predict phone-level human la-
076 beled scores. Compared with (Zhang et al., 2021),
077
078
079
080
081
082

the proposed system shows promising performance on the speechocean762¹ data set. A group of experiments also show better performance on Chinese children’s speech MD&D task than well tuned on-line Chain model based engine with smaller model size.

The remainder of this paper is organized as follows: Section 2 introduces the proposed Conformer based CAPT system. Section 3 presents the experiments and results. The conclusion is drawn in the Section 4.

2 Proposed CAPT Framework

Predominant approaches to CAPT primarily performs MD&D by extending ASR technologies, especially through post processing recognition scores. The proposed CAPT framework is shown in Figure. 1 which consists of acoustic model, feature extraction module and feedback module.

2.1 Acoustic Model

For the CAPT, acoustic models (e.g. DNN-HMM) are applied to infer acoustic features. Conformer (Gulati et al., 2020) based model draw immense interest recently and became the dominated model due to its ability to pay attention on both local and global dependencies of the utterance. The encoder is stacked by several blocks which consist of a positionwise feed-forward module, a multihead self-attention module, a convolution module, and another feed-forward module in the end. Only encoder architecture is used in our model to predict phoneme result and evaluate pronunciation for the latency-accuracy trade-off.

In order to achieve phoneme level MD task, the Conformer based encoder is employed to train phonetic level targets. It is necessary to convert transcriptions from word level to phoneme level. The method for obtaining phonetic representation of an input utterance is mainly based on the pronunciation dictionary. There are some notices in the conversion. If a word is not in the vocabulary, the out of vocabulary (OOV), the phonetic sequence of this word will be converted to the ’<unk>’ unit which represents unknown phones. It is also inevitable that there are some words with multiple pronunciation. In the conversion, the pronunciation will be decided by the part of speech which is implemented by HMM model (Hajic et al., 2009).

¹<https://www.openslr.org/101/>

2.2 Feature Extraction

In our Conformer-based system, we employ CTC-based end-to-end force alignment approach to compute phoneme level features. The log phone posterior ratio between the canonical phone and the one with the highest score is used to approximate GOP score as shown,

$$\text{GOP}(p) \approx \log \frac{P(p|\mathbf{o}; t_s, t_e)}{\max_{q \in Q} P(q|\mathbf{o}; t_s, t_e)} \quad (1)$$

where t_s and t_e are the start and end frame indexes, respectively; Q is the whole phone set; $P(p)$ is the prior of phone p ; \mathbf{o} is the acoustic aligned segment.

For the phoneme level scoring, LPP and LPR features are extracted by the output and the force alignments. The LPP of phone p is defined as

$$\text{LPP}(p) \approx \sum_{t=t_s}^{t_e} \log P(s_t|\mathbf{o}_t) / \text{NF}(p), \quad (2)$$

where $\text{NF}(p)$ is the number of frames in the phoneme p ; \mathbf{o}_t is the augmented input observations of the frame t ; s_t is the senone label of the frame t generated by force alignment with the given canonical phone p ; The LPR between phone p_i and p_j is defined as:

$$\text{LPR}(p_j|p_i) = \log P(p_j|\mathbf{o}; t_s, t_e) - \log P(p_i|\mathbf{o}; t_s, t_e). \quad (3)$$

For the sentence level scoring, some features are extracted from the LPP features. There are two feature groups coming from GOP related value and phoneme level recognized results by CTC decoder. The mean values and variances of GOP numerator, GOP denominator and GOP values are computed with the force aligned phoneme sequence. In addition to this, the recognized details including phoneme error rate (PER), the insertion PER, the deletion PER and the substitution PER comparing with the referred phoneme are used for scoring.

2.3 Feedback Module

For phoneme and sentence level scoring, a tree boosting system named XGBoost (Chen and Guestrin, 2016) is trained with the expert labeled score. The phoneme scores are used to detect the mispronunciation. The sentence score is used to represent the overall pronunciation evaluation.

For the phoneme level pronunciation diagnosis, the CTC decoder is applied CTC prefix beam search on the CTC output of the model, which can

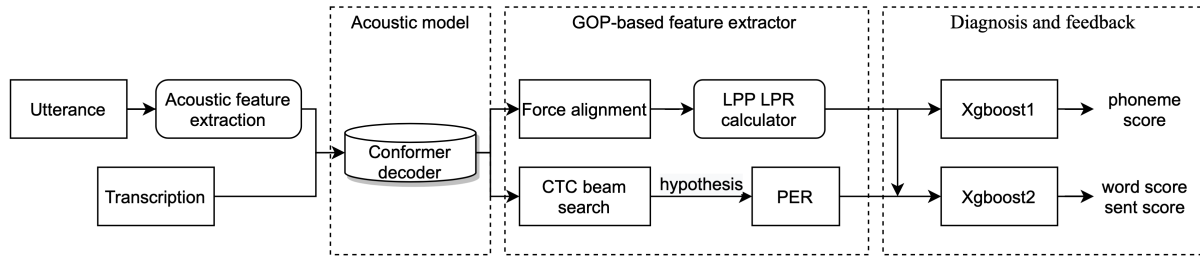


Figure 1: Framework of the proposed CAPT system

173 give the 1-best results. This phoneme level results
 174 are used for the diagnosis based on the reference
 175 phoneme sequence.

176 3 Experiment

177 3.1 Databases

178 With reproducibility in mind, open datasets are
 179 adopted. Librispeech (Panayotov et al., 2015) is
 180 used to train the native AM. A new open-source
 181 non-native english speech corpus named "spee-
 182 chocean762" (Zhang et al., 2021), recorded by L2
 183 language learners, is used for pronunciation assess-
 184 ment. To demonstrate the performance on non-
 185 native children speech, our internal dataset Yiqi
 186 Voice and Yiqi Evaluation are also included. The
 187 Yiqi Voice contains 3,000 hours Chinese children’s
 188 speech of designated transcription and conversation
 189 which are used to train AM. The Yiqi Evaluation
 190 contains about 100,000 utterances uniformly rang-
 191 ing from 0 to 8 in score labeled by experts which
 192 is used for MD&D. All of our internal data are
 193 anonymized and eligible to be used for research
 194 purposes.

195 3.2 Implementation Details

196 In the baseline model, 40 high-resolution mel-
 197 frequency cepstral coefficients and 100 dimen-
 198 sional i-vector features are extracted. We use 5
 199 layer time-delay neural networks (TDNN) with
 200 1280 dimensions. The AM output is used for the
 201 forced alignment and the computation to obtain
 202 the GOP values and the GOP-based features. As
 203 regards the evaluation model, support vector ma-
 204 chines (SVM) classifiers are built for each phone
 205 with the GOP-based features and the corresponding
 206 manual scores. Furthermore, chain model (Povey
 207 et al., 2016) is applied to compare with Conformer
 208 model. A 4-gram language model is trained for the
 209 word and phoneme level decoding.

210 In this paper, Wenet² is chosen to train byte
 211 pair encoding (BPE) and phoneme level Conformer
 212 models. 80 dimensional log-mel filterbank com-
 213 puted are extracted for the model input. In front of
 214 the encoder, two convolution sub-sampling layers
 215 are used with 4 times sub-sampling in total. We
 216 use 12 Conformer blocks in which have 4 multi-
 217 head attention with 64 output dimensions. For the
 218 XGBoost training, the learning rate is 0.1 and the
 219 weights of L1 regularization term is 0.001. The
 220 maximum tree depth is 7.

221 3.3 Recognition Performance

222 To evaluate the pronunciation, it is necessary to
 223 predict the phoneme level score together with the
 224 recognized phoneme level results for MD&D.

	Model Arch.	test-clean	test-other
WER	HMM-TDNN	4.3	10.62
	Conformer Attention Rescore	3.12	8.55
PER	HMM-TDNN	9.55	19.86
	Conformer Attention Rescore	1.84	5.48
	Conformer CTC Beam Search	1.87	5.61

Table 1: Comparison of different models in word and phoneme level.

	Model Arch.	Yiqi Voice
WER	Chain model	9.59
	Conformer	8.6

Table 2: Comparisons on Yiqi Voice dataset.

225 Table 1 and Table 2 show the comparisons of
 226 different models on word level and phoneme level.
 227 It can be seen that Conformer shows much better
 228 word error rate (WER) and PER than HMM-TDNN

²<https://github.com/mobvoi/wenet>

on Librispeech and Yiqi Voice. The PER performance is adequate to detect phoneme level pronunciation error and diagnose the phoneme details. In the Conformer decoding, the CTC beam search decoding method apply CTC prefix beam search on the CTC output and the attention rescore method rescoring the n-best candidates applied by the CTC beam search introduced in (Yao et al., 2021). It can be shown that the attention rescore performs better than CTC beam search method slightly. To balance the latency-accuracy trade-off, only the encoder part of AED is applied in this paper.

3.4 Mispronunciation Detection and Scoring Performance

The performance of the proposed system is compared with Kaldi recipe on the "speechocean762" dataset and our online CAPT system on the Yiqi Evaluation dataset.

3.4.1 Results on the Speechocean762

For evaluating the system's performance, the scoring results are shown in Table 3 comparing the baseline which can be found in Kaldi recipe³ (Zhang et al., 2021). It can be seen that our proposed system outperforms the baseline from 0.45 to 0.5 on Pearson correlation coefficient (PCC) and from 0.16 to 0.11 on mean squared error (MSE). Furthermore, the performance on sentence level on the "speechocean762" test sets shows 0.66 on PCC metrics.

	Level	MSE	PCC
HMM-TDNN	Phoneme	0.16	0.45
Proposed		0.11	0.50
Proposed	Sentence	1.399	0.654

Table 3: Comparisons of the evaluation performance between Kaldi recipe and our proposed system.

3.4.2 Results on the Yiqi Evaluation

To show the advantage of Conformer model in CAPT system, the performance of the proposed system is compared with our online well optimized system which consists of chain model tuned on Yiqi Voice, GOP module and XGBoost scoring model. The engine is implemented using fixed-point, with server-based multi-thread parallel method, for hundreds of millions MD&D requests per day. As

³https://github.com/kaldi-asr/kaldi/tree/master/egs/gop_speechocean762

shown in the Table 4, the Conformer performs better than the well optimized system with smaller model size.

Model Arch.	Model Size	MSE	PCC
Chain model	33M	0.944	0.90
Proposed	19M	0.802	0.91

Table 4: Comparisons on Yiqi Evaluation dataset

3.5 Diagnosis Advantage

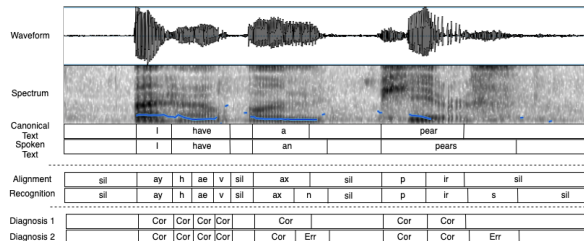


Figure 2: An example to detect and diagnose the insertion errors.

In the traditional alignment based MD&D system, some likely error patterns can be covered making use of the context-dependent phonological rules in the diagnosis. There are still some errors can not be diagnosed if they are not designed in the alignment graph. For example shown in Figure 2, when the reference is "I have a pear", the student pronounces "I have an pears". This kind of insertion error can not be detected and diagnosed based on the alignment methods (Diagnosis 1). In the proposed recognition based diagnosis (Diagnosis 2), the insertion errors like phoneme '/n/' and '/s/' can be detected and used to show the diagnosis for the learners.

4 Conclusion

In this paper, we proposed a Conformer/XGBoost based CAPT system providing MD&D. In the system, a phoneme level Conformer is trained with standard English corpus in this paper. After that, LPP, LPR, and some other features are extracted considering the reference phoneme sequence. Finally, we employ XGBoost to predict the pronunciation scores. The experiment results show that our proposed system outperforms the baseline system and well tuned online Chain model based system with smaller model size. In future, we will introduce works for the abundant dimension evaluations including scores of phoneme, word and sentence level, stress, fluency, prosody and so on.

References

- 300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
- Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794.
- Yiqing Feng, Guanyu Fu, Qingcai Chen, and Kai Chen. 2020. Sed-mdd: Towards sentence dependent end-to-end mispronunciation detection and diagnosis. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3492–3496. IEEE.
- Kaiqi Fu, Jones Lin, Dengfeng Ke, Yanlu Xie, Jinsong Zhang, and Binghuai Lin. 2021. A full text-dependent end to end mispronunciation detection and diagnosis with easy data augmentation techniques. *arXiv preprint arXiv:2104.08428*.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*.
- Jan Hajic, Jan Raab, Miroslav Spousta, et al. 2009. Semi-supervised training for the averaged perceptron pos tagger. In *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*, pages 763–771.
- Wenping Hu, Yao Qian, Frank K Soong, and Yong Wang. 2015. Improved mispronunciation detection with deep neural network trained acoustic models and transfer learning based logistic regression classifiers. *Speech Communication*, 67:154–166.
- Shao-Wei Fan Jiang, Bi-Cheng Yan, Tien-Hong Lo, Fu-An Chao, and Berlin Chen. 2021. Towards robust mispronunciation detection and diagnosis for 12 english learners with accent-modulating methods. *arXiv preprint arXiv:2108.11627*.
- Wai-Kim Leung, Xunying Liu, and Helen Meng. 2019. Cnn-rnn-ctc based end-to-end mispronunciation detection and diagnosis. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8132–8136. IEEE.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE.
- Patrizia Paterlini-Brechot and Naoual Linda Benali. 2007. Circulating tumor cells (ctc) detection: clinical impact and future directions. *Cancer letters*, 253(2):180–204.
- Daniel Povey, Vijayaditya Peddinti, Daniel Galvez, Pegah Ghahremani, Vimal Manohar, Xingyu Na, Yiming Wang, and Sanjeev Khudanpur. 2016. Purely sequence-trained neural networks for asr based on lattice-free mmi. In *Interspeech*, pages 2751–2755.
- Sweekar Sudhakara, Manoj Kumar Ramanathi, Chiranjeevi Yarra, and Prasanta Kumar Ghosh. 2019. An improved goodness of pronunciation (gop) measure for pronunciation evaluation with dnn-hmm system considering hmm transition probabilities. In *INTER-SPEECH*, pages 954–958.
- Zhenyu Wana, John HL Hansen, and Yanlu Xie. 2020. A multi-view approach for mandarin non-native mispronunciation verification. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8079–8083. IEEE.
- Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R Hershey, and Tomoki Hayashi. 2017. Hybrid ctc/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253.
- Silke M Witt and Steve J Young. 2000. Phone-level pronunciation scoring and assessment for interactive language learning. *Speech communication*, 30(2-3):95–108.
- Silke Maren Witt. 1999. Use of speech recognition in computer-assisted language learning.
- Bi-Cheng Yan, Meng-Che Wu, Hsiao-Tsung Hung, and Berlin Chen. 2020. An end-to-end mispronunciation detection system for 12 english speech leveraging novel anti-phone modeling. *arXiv preprint arXiv:2005.11950*.
- Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang, Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu Chen, Lei Xie, and Xin Lei. 2021. Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit. *arXiv e-prints*, pages arXiv-2102.
- Junbo Zhang, Zhiwen Zhang, Yongqing Wang, Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li, Daniel Povey, and Yujun Wang. 2021. speechocean762: An open-source non-native english speech corpus for pronunciation assessment. *arXiv preprint arXiv:2104.01378*.