

# Multilingual Speech Recognition Based on The End-to-End Framework

Sendong Liang · Wei Qi Yan

Received: 01 June 2021 / Accepted: 01 June 2021

**Abstract** Speech recognition is an important field in natural language processing. In this paper, the end-to-end framework for speech recognition with multilingual datasets is proposed. The focus of this paper is on the end-to-end methods. Our objective is to improve the performance of the CTC+Attention model. In order to compare speech recognition methods for multiple languages, we design and create three datasets: Chinese, English and Code-Switch. We evaluate the proposed hybrid CTC+Attention model in multilingual environment. Throughout our experiments, we find that the end-to-end framework of the hybrid CTC+Attention model achieves better performance compared to the HMM-DNN model in a single language and Code-Switch speaking environment.

**Keywords** Speech recognition · End-to-end framework · Attention model · CTC model · Code-Switch

## 1 Introduction

### 1.1 Background

Spoken language is essential to modern human cultures. Speech is a way of communications amongst people through languages. In the past decades, with the development of intelligent devices, the developed applications began to enter people's daily life. When people communicate with intelligent devices through languages, automatic speech recognition(ASR) plays the pivotal role. ASR aims to convert the audio content into corresponding text, the text is further processed through human-computer interaction, such as multilingual translations and hand talk, etc.

In the early stage, because it was impossible to directly model the audio-to-text conversion, Bayes' theorem was implemented to convert human speech into text

---

Sendong Liang  
E-mail: ghg0412@autuni.ac.nz

Wei Qi Yan  
E-mail: weiqi.yan@aut.ac.nz

so as to calculate the corresponding audio features. Accordingly, the probability of audio feature sequences is decomposed to the product of conditional probabilities of corresponding audio features.

With the development of HMM, speech recognition has transited from isolated words of a small system for speech recognition to a large vocabulary continuous system nowadays [37]. The framework based on HMM model for speech recognition shows its excellence and reliable stability, which was the mainstream speech recognition model.

Deep neural network has contributed to acoustic models and formed the HMM-DNN speech recognition framework [18]. The objective function of the DNN-based acoustic model is the probability of a HMM state given a sequence of audio features. With more and more research work dedicated to deep learning models, the ability of acoustic models becomes stronger and stronger. A classification network has been directly created from the HMM-DNN acoustic model [8].

Owing to the soaring deep learning, the end-to-end frameworks for speech recognition have shown exemplary performance in high-resource languages. However, it is hard to perform well in low-resource datasets for speech recognition, such as Chinese-English and Code-Switch environment [38]. In this paper, we will devote to bridge the gap and provide our solution for resolving this problem.

## 1.2 Motivation

Up to date, our daily communications often are surrounded by mixed languages, which is academically named as Code-Switch model. For example, a lot of Chinese people mingle English words in a Chinese sentence. The speech blended with the words from multilingual is one of the critical challenges in speech recognition. The main technical difficulty also includes the non-native accents, the composition brings difficulties to modelling the mixed acoustics, meanwhile, the labelled datasets for the mixed speech recognition are exceptionally scarce.

Traditional phonetic framework is based on basic acoustic unit for language recognition. The linguistic information is different for different languages, such as Chinese phonetic consonants and English phonemes. The framework relies on the specific linguistic knowledge and is tough to be expanded to multilingual speech recognition. The end-to-end frameworks employ a unified network for modelling, require less manual dictionary editing, and are dependent more on datasets than linguistic information. Accordingly, we have a great interest in using the end-to-end framework to solve this emerging speech recognition issues.

The remaining part of this paper is organized as follows. We have our literature review in Section 2, our method is present in Section 3, our result analysis is shown in Section 4, our conclusion is drawn in Section 5.

## 2 Literature Review

### 2.1 Automatic Speech Recognition

Speech recognition models are grouped into three categories. The first category includes the rule-based models, such as Shoebox and Harpy created by IBM and

CMU in 1962 and 1976, respectively [11]. The second group encompasses the statistical-based models such as Large Vocabulary Continuous Speech Recognition (LVCSR) and HMM. HMM-Gaussian Mixture Model (GMM) has been the dominant framework in the field of speech recognition till the emerge of deep learning [34]. The third one is deep-learning-based models such as Deep Neural Networks (DNNs), Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), and bidirectional LSTM (BiLSTM). LSTM is sensitive to the static data, which may lead to delays with respect to features, BiLSTM appears as a special architecture operating the input sequence in both directions [17].

In 2006, an unsupervised method was employed to pretrain the Deep Belief Networks (DBNs), which solved the problem that gradient descent is sensitive to the initial value [9]. DNNs have been applied to acoustic modelling to form the HMM-DNN framework in speech recognition [18]. The objective function of DNN-based acoustic models is the probability of a HMM state given a sequence of audio features. In 2009, DNNs were applied to the TIMIT phoneme recognition and achieved excellent performance [23]. In 2012, CNNs were applied to the LVCSR, which normalized the data and obtained a higher performance in speech recognition [1].

Moreover, an end-to-end framework for speech recognition based on RNNs and Weighted Finite State Transducer (WFST) decoding method was proposed in 2015. This end-to-end framework directly utilises analogue signals as its input, which improves the recognition rate [21] instantaneously. In 2017, a new method of end-to-end speech recognition was proffered, which employs CTC model in a multitask framework to improve the robustness of the system and achieves quick convergence, therefore alleviates the problem of data alignment [12].

However, most of the previous research work in speech recognition focuses on better performance of speech recognition for a single language or a single task. The comparison of speech recognition between multiple languages, the publications of speech recognition in a multilingual or Code-Switch environment are relatively rare.

## 2.2 Language Model Using RNNs

Feedforward Neural Networks (FFNNs) are dependent on the preceding words, which cannot learn the dependent information of long sentences. With the development of deep learning, RNNs turn up, which put forward the issue of long sequence dependence. RNNs are defined as a type of ANNs that make classification and predictions for various data, such as text, audio, video, genomes, etc. [19].

Compared with traditional FFNNs, such as multilayer perception, using static classifiers and only considering fixed-size input windows which are irrespective of surrounding context, RNNs are more effective and suitable to transcribe time series such as speech transcription because of its hidden network layers [6]. Different from FFNNs using fixed-length context, RNNs do not take use of a limited size of context, which contain cache models that encode temporal information implicitly with arbitrary lengths [22]. The recurrent connections allow information to cycle inside networks for a long time adapting to the past inputs [2]. A continuous

vector space best suits for word representation, meanwhile, deep learning methods reflect the relationship between continuous words. Accordingly, compared to FFNNs, RNNs best solve the contextual dependency problem which spans over a fixed number of predecessor words [30].

## 2.3 The End-to-End Framework for Speech Recognition

The end-to-end framework for speech recognition refers to directly transduce the input sequence of acoustic feature vectors to the output sequence of token such as phonemes, characters, or words [14]. The end-to-end model is split into three categories based on the alignment methods: Connectionist Temporal Classification (CTC), Attention Encoder-Decoder (AED), and RNN Transducer (RNN-T), which have been widely utilized in large-scale speech recognition [14]. The end-to-end models are more suitable for on-device applications than conventional speech recognition because there are fewer parameters by folding the acoustics, pronunciation, and language models into one neural network [13].

## 2.4 Attention Encoder-Decoder Model

Attention-based encoder-decoder model such as LAS (Listen, Attend and Spell) contains three main components: Encoder as an acoustic model, attender as an alignment model, and decoder as a language model [4], which subsume the acoustics, pronunciation and language model components into a single neural network without a lexicon or a separate text normalization component [5].

Applied to Google voice search, the proposed model achieves a WER of 5.6%, while the hybrid HMM-LSTM model attains 6.7% WER. Throughout testing the same models based on dictation, the proposed model reaches up to 4.1%, the HMM-LSTM model gets 5% WER. The decoding process of sequence-to-sequence (S2S) models with soft attention incurs a quadratic time and space cost, which is regarded as a challenge for online sequence transduction [5].

In order to address the online streaming challenge of the attention-based model, monotonic-chunkwise attention was put forward, which splits input sequences into numbers of small chunks [5]. Triggered attention equipped with the CTC-based classifier performs well to control the activation of the attention-based decoder [24].

# 3 Methods

## 3.1 Data Preparation

### 3.1.1 Language Features

By considering English and Mandarin are worldwide widely used languages in the low-resource environment of bilingual Code-Switch corpus, it is important to select appropriate recognition units for acoustic modelling, which convert the

feature vector sequence to the speech recognition unit by recognizing the speech recognition unit corresponding to the feature vector sequence [29].

English is an Indo-European language, while Mandarin is a Sino-Tibetan language [3]. Based on the Oxford Dictionary, English is written in Latin alphabets, namely, the Roman alphabets, which contain 26 letters and nearly 170,000 words. The general modelling units in English include phone, subword, and character [35].

In the field of Chinese speech recognition, there are various available acoustic modelling units, including Chinese characters (word), syllable (syllable), semi-syllable (initial/final), phoneme (phone), which is generally based on phonetic knowledge or data-driven generation [41]. Mandarin contains more than 6,000 characters, 60 phonemes, 408 atonal syllables, and 1,302 toned syllables [16]. Each syllable includes initials, finals and tones, in total, Chinese Mandarin 22 initials and 39 finals of syllables.

Meanwhile, there are a plenty of homophones and polyphones in Chinese [42], which may need high-level non-acoustic context knowledge for speech recognition. A small flexible unit may lead to the difficulties of calibrating the training dataset. By contrast, the limitation of flexibility as well as the high requirement of lexicon and out of vocabulary (OOV) [32] are significant issues though the large unit equipped with high recognition performance. The syllable unit meets the requirement and flexibility [28]. Moreover, the utilization of syllables with tones effectively increases the recognition accuracy compared to Chinese characters and syllable initial/final with tones [7].

### 3.2 Corpus Design

Based on the language features in Section 3.1.1, three speech datasets are applied to our experiments, which include dataset Alpha (Mandarin), dataset Beta (English), and dataset Gamma (Mandarin-English).

We create the three datasets through the text-to-speech iFLYTEK InterPhonic toolkit. It is a software developed by iFLYTEK[10], which converts text into male or female voices. The toolkit is based on an advanced large corpus and a phonetic prosody description system, the quality of the synthesized voices in *.wav* format files is comparable to that of a real person.

Converting the transcript text to synthesised acoustic speeches controls the variables. For example, speaker's accents, speaker's emotions, and environmental noises are consistent. The focus of our experiments by using the synthesised acoustic audio is on the recognition results of multilingual environment without considering useless variables or factors.

#### 3.2.1 Dataset Alpha (Mandarin)

Alpha is a Mandarin speech dataset. We created this dataset through the iFlytek toolkit based on the transcripts of the THCHS-30 corpus. In the THCHS-30 corpus, there are numerous transcript files corresponding to the acoustic files. Each file contains three rows, which is the correctly-labeled transcript of the corresponding acoustic sentence. The first row is the Chinese characters of the corresponding

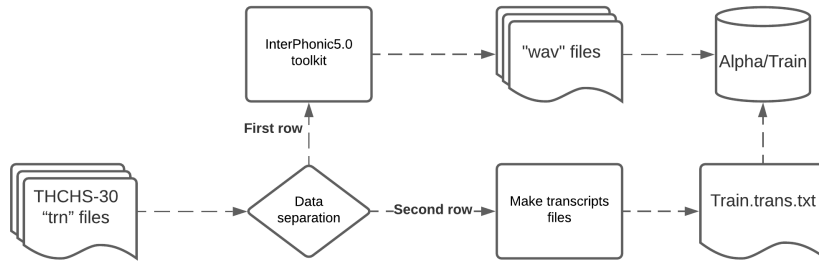


Fig. 1: The process of creating Alpha-Train

acoustic sentence. The second row is the Chinese Pinyin with a tone, which corresponds to the Chinese characters. The third row is the syllable initial-final tone.

As shown in Fig. 1, we selected 8,000 files of the THCHS-30 corpus. The first row (Chinese characters) of each selected file are collected as the input of the iFlytek InterPhonic toolkit. The selected Chinese characters are synthesised into 8,000 files sampled at 16 kHz. The second row (Chinese Pinyin with tones) of each selected file is collected together as one text file as our labeled transcript for training. Throughout this operation process, we created the training subset based on Alpha (Mandarin).

Similarly, we created the “Dev” (development) subset and the “Test” subset from Alpha (Mandarin) through the same process. The “Dev” subset is synthesised from other 1,000 files of the THCHS-30 corpus except for the selected 8,000 files. The “Test” subset is synthesised from other 1,000 transcript files except for the selected 9,000 files.

The synthesised acoustic files in the dataset Alpha (i.e., Mandarin) are divided into three groups as shown in Table 1: One-hour acoustic files as the test subset (i.e., Alpha-Test); One-hour acoustic files as the development subset (Alpha-Dev) to train the rapid reaction and performance evaluation, the other 10-hour acoustic files as the training subset (i.e., Alpha-Train).

Table 1: The statistics of dataset Alpha (Mandarin)

Datasets	Time(Hours)	Speakers	Male	Female	Percentages of THCHS-30
Alpha-Train	10	2	1	1	29.878%
Alpha-Dev	1	2	1	1	2.988%
Alpha-Test	1	2	1	1	2.988%

### 3.2.2 Dataset Beta (English)

Dataset Beta (English) is an English speech dataset. We created this dataset by using the iFlytek InterPhonic toolkit based on the transcript files of the LibriSpeech corpus [25]. The production method is similar to Alpha as shown in Fig. 2. The structure of the Beta dataset is shown in Table 2: 10-hour acoustic files synthesised

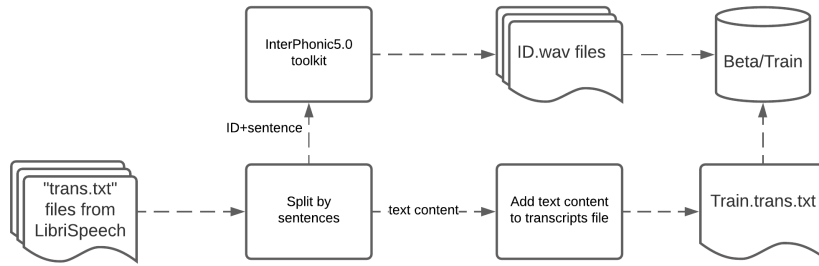


Fig. 2: The process of creating dataset Beta-Train

as the training subset (i.e., Beta-Train), one-hour audio files as the development subset (i.e., Beta-Dev) to train the rapid reaction and performance evaluation, and one-hour audio files as the testing subset (i.e., Beta-Test).

Table 2: The dataset Beta (English)

Dataset	Time(hour)	Speakers	Male	Female	Percentage of LibriSpeech-clean
Train	10	2	1	1	8.977%
Test	1	2	1	1	0.898%
Dev	1	2	1	1	0.898%

### 3.2.3 Dataset Gamma (Mandarin-English)

Dataset Gamma (Mandarin-English) is a mixed Mandarin-English bilingual speech database. We created this dataset through the iFlytek InterPhonic toolkit based on the transcript files of the TAL-CSASR corpus [31]. The original TAL-CSASR corpus is the audio captured in English class teaching environment. The production method and folder structure of the Gamma dataset are similar to Alpha and Beta. The production method is shown in Fig. 3. The folder structure of the Gamma dataset is as same as Alpha and Beta. The details of each sub-dataset are shown in Table 3.

Table 3: The dataset Gamma (Mandarin-English)

Datasets	Time(hours)	Speakers	Male	Female
Train	10	2	1	1
Dev	1	2	1	1
Test	1	2	1	1

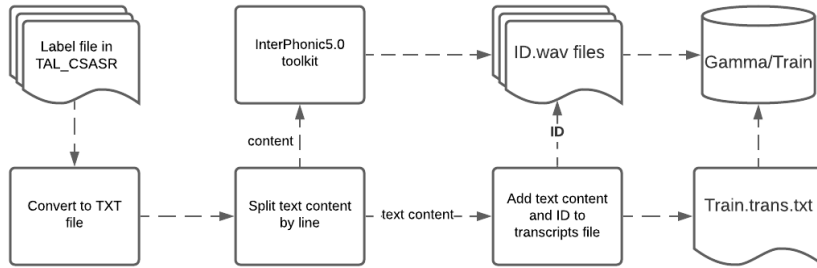


Fig. 3: The process of creating dataset Gamma

### 3.3 Experiment

#### 3.3.1 Experiment Environment and Setup

We implement Kaldi [27] in the experiment of traditional speech recognition. The end to end framework takes use of the PyTorch and ESPnet [36] frameworks based on LAS[4] in the RNN language model. Based on the inspiration [33], we superimposed multihed and location attention in the hybrid framework. Table 4 shows the details of the hardware configuration and operating system.

Table 4: The hardware Equipment of Experiments

OS Name	Ubuntu20.04
OS Type	64-bit
Memory	32 GiB
Processor	Inter(R)Core i9-9900k CPU @ 3.60GHZ x 16
GPU	GeForce RTX 2080Ti 11GiB x 2
Disk Capacity	1.5T

#### 3.3.2 Evaluation Methods

The most straightforward approach to evaluate the performance of speech recognition is to calculate the word error rate (WER) [20]. WER is employed for the English dataset, which is computed through Eq. (1).

$$R = \frac{I_E + D_E + S_E}{N} \times 100\% \quad (1)$$

where  $I_E$  refers to the insertion number of English words,  $D_E$  means the deletion number of English words,  $S_E$  stands for the substitution number of English words, and  $N$  is the total number of English words in the correct sentence. However, characters are considered instead of words in the Chinese THCHS-30 corpus. Similarly, the character error rate (CER) is calculated through Eq. (2), in which the suffixes



ended with  $M$  take the place of  $E$ .

$$R = \frac{I_M + D_M + S_M}{N} \times 100\% \quad (2)$$

### 3.3.3 Parameter Optimization

The main idea of the CTC+Attention hybrid model is to utilise CTC to force the alignment of the eigenvectors of the audio frames to reduce a single tag [26]. At the same time, CTC does not allow skipping label output under the same audio characteristics to avoid the frame skipping mentioned in the attention subsection. The score function of the hybrid model is described by using Eq. (3) based on the schematic diagram of the combination of CTC and attention formulas. Among them,  $O_{hybrid}$  is the prediction result of this model,  $Y$  is the text label sequence,  $X$  is the feature vector sequence corresponding to the audio frame,  $\lambda$  is the evaluation weight of the CTC model.  $\log P_c(Y | X)$  is the score function of CTC, and  $\log P_a(Y | X)$  is the score function of the attention model. Therefore, the optimization process is regarded as finding the optimal solution of  $\lambda$ .

$$O_{hybrid} = \log P_a(Y | X)(1 - \lambda) + \log P_c(Y | X)\lambda \quad (3)$$

Five model trainings were carried out for  $\lambda$  from 0.2 to 0.7 based on Alpha dataset. The losses with different  $\lambda$  are shown in Fig. 4. The red line is the loss of the attention model, the blue line is the loss of CTC. Judged from the loss of CTC, the convergence speed of the five experiments has been improved. While CTC improves the learning efficiency, the convergence of attention does show relatively large fluctuations. If  $\lambda$  equals to 0.7, the loss of attention has a huge fluctuation in a short period.

On the other hand, the CERs of the five experiments based on the Alpha-Dev subset are listed in Table 5. If  $\lambda$  is less than 0.6, the CER obtained with the increase of  $\lambda$  becomes smaller, from 8.85% to 8.32%. However, if  $\lambda$  is equal to 0.6, CER starts increasing. Unlike the slight decrease in attention model, the performance of CTC regarding CER dropped from 14.28% to 11.04% as  $\lambda$  on the raise. The minimum is reached if  $\lambda$  is equal to 0.6. Although  $\lambda$  is equal to 0.5, attention model has achieved the best results. Nevertheless, compared to the slight difference in attention, if  $\lambda$  is equal to 0.6, the CTC improvement is much remarkable. Combined the experimental results in Fig. 4 and Table 5, we see that  $\lambda$  equals 0.6 in the mixed model is the optimal solution.

Table 5: The comparisons of CERs based on ‘‘Dev’’ subset with various CTC weights

$\lambda$	CER(Attention)	CER(CTC)
0.2	8.85%	14.28%
0.4	8.52%	12.84%
0.5	<b>8.24%</b>	12.02%
0.6	8.32%	<b>11.04%</b>
0.7	8.51%	11.17%

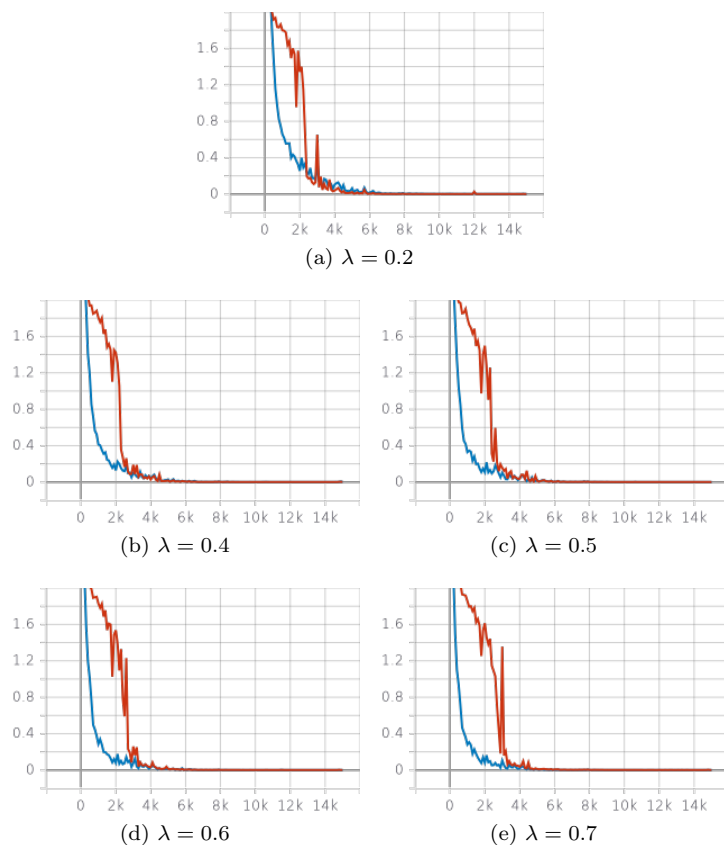


Fig. 4: The comparisons of loss with various CTC weights

### 3.4 Experimental Results

The evaluations were carried out based on three datasets: Alpha, Beta, and Gamma, respectively. The final experimental results and the previous three experimental results are compared in Table 6.

Table 6: The comparisons of experimental results(CERs/WERs)

Models	Alpha	Beta	Gamma
HMM-TDNN-F	10.91%	<b>18.23%</b>	<b>25.62%</b>
CTC+Attention	<b>10.22%</b>	19.05%	26.11%

## 4 Result Analysis

### 4.1 Comparisons of Speech Recognition Frameworks

The first row in Table 6 shows that the traditional HMM-based speech recognition model achieved 10.91% error rate in Chinese and 18.23% in English. These two items are much higher than only implementing CTC or attention model in speech recognition. However, the CTC+Attention hybrid model obtained better results based on the Chinese dataset, with the error rate 0.069% less than that of the traditional model. Based on the Code-Switch dataset, the result of the traditional speech recognition model regarding word error is the best of the three models, but the difference with CTC+Attention is not obvious, and the outcome is improved by 3%. The evaluation results of all models in the Gamma dataset are worse than Alpha and Beta. That is due to the complexity of mixed languages. Moreover, through the experimental results, we see that the accuracy of CTC+Attention hybrid model is the highest one, an error rate 10.22% was achieved based on the Chinese dataset.

### 4.2 Recognition Rate of Multiple Languages

From the language perspective, the experimental results show that the English dataset is more challenging to be used for model training than the Chinese one. Compared with the three languages, Chinese labels with Pinyin performs better than English in speech recognition. The convergence in the training phase and the WERs in the verification phase consistently reflect that the Chinese dataset Alpha labeled with Pinyin is easier to be trained, the recognition accuracy is greater than English. The utilization of Pinyin effectively avoids the situation of one sound and multiple characters in Chinese.

On the other hand, the WER of the HMM-TDNN-F model in the dataset Gamma is 25.62%. In contrast, the CTC+Attention model has a 37% gap between the two datasets. The results show that the difference in recognition accuracy of the CTC+Attention model by using different language datasets is smaller. The difference of multilingual speech recognition is better than the traditional framework.

## 5 Conclusion

Throughout this paper, we provide a research foundation for future exploration of speech recognition in different languages. The comparison of actual outcomes in multiple languages provides a future research direction based on the characteristics of speech recognition.

In this paper, we investigate the speech recognition performance based on the CTC and attention hybrid model by using our three labeled datasets. The model attains 10.91%, 18.23%, and 25.62% error rates based on the Chinese, English, and Chinese-English Code-Switch datasets. The CTC+Attention model adopts the optimal solution to complete the model evaluations and achieves similar performance to the traditional model. The evaluation based on the Chinese dataset defeated

the traditional model and reached 10.22% CER. Although the performance based on the English and Code-Switch datasets was not as good as the traditional model, the gap remained within 3%.

In future, more attention model will be added to our project to replace the current model. We will implement the predetermined algorithms to improve the training results. In addition, we will explore multiple languages for speech recognition by using the methods of creating datasets in this paper [15,40,39].

## References

1. Abdel-Hamid, O., Mohamed, A.R., Jiang, H., Penn, G.: Applying convolutional neural networks concepts to hybrid NN-HMM model for speech recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4277–4280. IEEE (2012)
2. Boden, M.: A guide to recurrent neural networks and backpropagation. The Dallas Project: SICS Technical Report (2002)
3. Chan, J.Y., Ching, P., Lee, T., Meng, H.M.: Detection of language boundary in code-switching utterances by bi-phone probabilities. In: International Symposium on Chinese Spoken Language Processing, pp. 293–296. IEEE (2004)
4. Chan, W., Jaitly, N., Le, Q.V., Vinyals, O.: Listen, attend and spell. arXiv preprint arXiv:1508.01211 (2015)
5. Chiu, C.C., Raffel, C.: Monotonic chunkwise attention. arXiv preprint arXiv:1712.05382 (2017)
6. Eyben, F., Wöllmer, M., Schuller, B., Graves, A.: From speech to letters-using a novel neural network architecture for grapheme based ASR. In: IEEE Workshop on Automatic Speech Recognition & Understanding, pp. 376–380. IEEE (2009)
7. Fu, L., Li, X., Zi, L.: Research on modeling units of transformer transducer for Mandarin speech recognition. arXiv preprint arXiv:2004.13522 (2020)
8. Georgescu, A.L., Cucu, H., Burileanu, C.: Kaldi-based DNN architectures for speech recognition in Romanian. In: International Conference on Speech Technology and Human-Computer Dialogue (SpeD), pp. 1–6. IEEE (2019)
9. Hinton, G.E., Osindero, S., Teh, Y.W.: A fast learning algorithm for deep belief nets. *Neural Computation* **18**(7), 1527–1554 (2006)
10. iFLYTEK Co., Ltd: Online TTS WebAPI. Website (2020). [https://global.xfyun.cn/products/online\\_tts](https://global.xfyun.cn/products/online_tts)
11. Jason, C.A., Kumar, S.: An appraisal on speech and emotion recognition technologies based on machine learning. *Language* **67**, 68 (2020)
12. Kim, S., Hori, T., Watanabe, S.: Joint CTC-attention based end-to-end speech recognition using multi-task learning. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 4835–4839. IEEE (2017)
13. Li, B., Chang, S.y., Sainath, T.N., Pang, R., He, Y., Strohman, T., Wu, Y.: Towards fast and accurate streaming end-to-end ASR. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 6069–6073. IEEE (2020)
14. Li, J., Zhao, R., Hu, H., Gong, Y.: Improving RNN transducer modeling for end-to-end speech recognition. In: IEEE Automatic Speech Recognition and Understanding Workshop, pp. 114–121. IEEE (2019)
15. Liang, S.: Multilingual Speech Recognition Based on the End-to-End Framework (Master’s Thesis). Auckland University of Technology, New Zealand (2021)
16. Lin, C.H., Lee, L.S., Ting, P.Y.: A new framework for recognition of Mandarin syllables with tones using sub-syllabic units. In: IEEE International Conference on Acoustics, Speech, and Signal Processing, vol. 2, pp. 227–230. IEEE (1993)
17. Liu, Z., Chen, Q., Hu, H., Tang, H., Zou, Y.: Teacher-student learning and post-processing for robust BiLSTM mask-based acoustic beamforming. In: International Conference on Neural Information Processing, pp. 522–533. Springer (2019)
18. Maas, A.L., Qi, P., Xie, Z., Hannun, A.Y., Lengerich, C.T., Jurafsky, D., Ng, A.Y.: Building DNN acoustic models for large vocabulary speech recognition. *Computer Speech & Language* **41**, 195–213 (2017)

19. Manaswi, N.K., Manaswi, N.K., John, S.: Deep Learning with Applications Using Python. Springer (2018)
20. Mansikkaniemi, A.: Acoustic Model and Language Model Adaptation for A Mobile Dictation Service (Master's Thesis). Aalto University (2010)
21. Miao, Y., Gowayed, M., Metze, F.: EESSEN: End-to-end speech recognition using deep RNN models and WFST-based decoding. In: IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 167–174. IEEE (2015)
22. Mikolov, T., Karafiát, M., Burget, L., Černocký, J., Khudanpur, S.: Recurrent neural network based language model. In: Annual Conference of the International Speech Communication Association (2010)
23. Mohamed, A.r., Dahl, G., Hinton, G.: Deep belief networks for phone recognition. In: NIPS Workshop on Deep Learning for Speech Recognition and Related Applications, vol. 1, p. 39. Vancouver, Canada (2009)
24. Moritz, N., Hori, T., Le Roux, J.: Triggered attention for end-to-end speech recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 5666–5670. IEEE (2019)
25. Panayotov, V., Chen, G., Povey, D., Khudanpur, S.: Librispeech: An ASR corpus based on public domain audio books. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 5206–5210. IEEE (2015)
26. Petridis, S., Stafylakis, T., Ma, P., Tzimiropoulos, G., Pantic, M.: Audio-visual speech recognition with a hybrid ctc/attention architecture. In: IEEE Spoken Language Technology Workshop, pp. 513–520. IEEE (2018)
27. Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., et al.: The Kaldi speech recognition toolkit. In: IEEE Workshop on Automatic Speech Recognition and Understanding. IEEE Signal Processing Society (2011)
28. Qu, Z., Haghani, P., Weinstein, E., Moreno, P.: Syllable-based acoustic modeling with CTC-SMBR-LSTM. In: IEEE Automatic Speech Recognition and Understanding Workshop, pp. 173–177. IEEE (2017)
29. Senior, A., Sak, H., Shafran, I.: Context dependent phone models for LSTM RNN acoustic modelling. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 4585–4589. IEEE (2015)
30. Sundermeyer, M., Oparin, I., Gauvain, J.L., Freiberger, B., Schlüter, R., Ney, H.: Comparison of feedforward and recurrent neural network language models. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 8430–8434. IEEE (2013)
31. TAL Education Group: TAL CS Auto Speech Recognition Data set. Website (2019). <https://ai.100tal.com/dataset>
32. Ueno, S., Inaguma, H., Mimura, M., Kawahara, T.: Acoustic-to-word attention-based model complemented with character-level CTC-based model. In: IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 5804–5808. IEEE (2018)
33. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. arXiv preprint arXiv:1706.03762 (2017)
34. Wang, D., Wang, X., Lv, S.: An overview of end-to-end automatic speech recognition. *Symmetry* **11**(8), 1018 (2019)
35. Wang, W., Wang, G., Bhatnagar, A., Zhou, Y., Xiong, C., Socher, R.: An investigation of phone-based subword units for end-to-end speech recognition. arXiv preprint arXiv:2004.04290 (2020)
36. Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Soplin, N.E.Y., Heymann, J., Wiesner, M., Chen, N., et al.: ESPnet: End-to-end speech processing toolkit. arXiv preprint arXiv:1804.00015 (2018)
37. Woodland, P.C., Odell, J.J., Valtchev, V., Young, S.J.: Large vocabulary continuous speech recognition using HTK. In: IEEE International Conference on Acoustics, Speech and Signal Processing, vol. 2, pp. II–125. IEEE (1994)
38. Wu, C.H., Shen, H.P., Yang, Y.T.: Chinese-English phone set construction for code-switching ASR using acoustic and DNN-extracted articulatory features. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **22**(4), 858–862 (2014)
39. Yan, W.: Introduction to Intelligent Surveillance Surveillance: Data Capture, Transmission, and Analytics. Springer (2019)
40. Yan, W.: Computational Methods for Deep Learning: Theoretic, Practice and Applications. Springer (2021)

41. Zenkel, T., Sanabria, R., Metze, F., Waibel, A.: Subword and crossword units for CTC acoustic models. arXiv preprint arXiv:1712.06855 (2017)
42. Zheng, Y., Yang, X., Dang, X.: Homophone-based label smoothing in end-to-end automatic speech recognition. arXiv preprint arXiv:2004.03437 (2020)