

# Improvement the phonetic recognition with sequence-length MFCC features and deep bidirectional LSTM

Toan Pham Van

Framgia Inc R&D Group

pham.van.toan@framgia.com

Hau Nguyen Thanh

Framgia Inc R&D Group

nguyen.thanh.hau@framgia.com

Thanh Ta Minh

Le Quy Don Technical University

ta.minh.thanh@framgia.com

**Abstract**—Phonetic recognition is an meaningful problem in the many fields of speech analysis. This applications can mention as dialect identification [1], mispronunciation detection [2], spoken document retrieval [3] etc. Some main approaches to solving this problem are improve the feature selection on input speech [4], apply some deep learning techniques [5][6][7] or combination both of them [8]. With the sequence data as the phonetics, the architecture based on recurrent neural network is an appropriate approach [9]. It is more powerful when combined with the improvement of features selection in input data. In our approach, we combining the Mel Frequency Cepstral Coefficients (MFCC) method with sequence-length for present the acoustic features of speech and using some RNN models to phonetic classification. All the experiments are implement on the Texas Instruments Massachusetts Institute of Technology (TIMIT) [10] phone recognition dataset. Especially, our data processing and features selection method give consistently better results with other researches when applied the same neural network model. Currently, we achieved the lowest error test rate - 15.24% by using Bidirectional LSTM, which is currently the best result in TIMIT dataset and reduction of about 1.3% over the last best result [5][6].

**Keywords**—Phonetic Recognition, MFCC features, sequence-length, bidirectional LSTM, TIMIT

## I. INTRODUCTION

Millions of people over the world study at least one foreign language. However, many of them can't approach to methods for mastering proper pronunciations. Nowadays, with the powerful of computer systems and internet, an automated computer tool to help language learners is really necessary. It try to detect mistakes made by non-native learners at either the phone, word, or sentence and inform the learner of those errors. Our approach is to detect phone-level mispronunciations in words detected by an Automatic Speech Recognition (ASR) system. The combination of CNN and RNN was constructed and and we chose a language which was not popular in previous researches is Japanese as an experiment for our solution.

## II. RELATED WORKS

### A. Baseline

A simple approach to solve this problem is using phoneme-based recognition and identifying pronunciation errors in the input speech of non-native speakers. But in actually, the accuracy for detection words is much higher than phonetic detection even for native speakers. It also mean that we can

not directly apply **Automatic Speech Recognition** system for mispronunciation detection. Instead we add to the ASR system a pronunciation model with possible faulty pronunciation variations are used to recognize the most likely phone sequences when it knew previous phones. Finally, the mispronunciation detection system is worked by executing the forced-alignment of ASR with extended pronunciation recognizer based on possible phonetic confusions. The figure below is a simulation of workflow in this system:

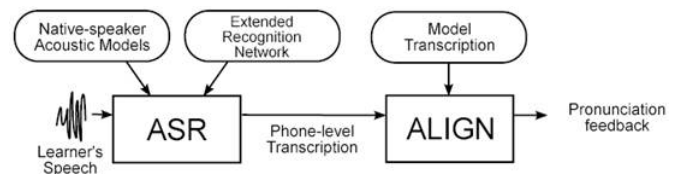


Fig. 1: ASR-based system to detect and diagnose L2 learners mispronunciation

### B. Feature Extraction

As discussed above, phone-level mispronunciation detection can detect mispronunciations in units of phones, words or sentences. Firstly, the speaker's speech samples are first converted to certain types of features such as Linear predictive cepstral coefficients (LPCC), Mel-frequency cepstral coefficients (MFCC), Power spectral analysis (FFT) or Mel scale cepstral analysis (MEL) etc. These features are used as the input of classifier. These allow the system to classify the mispronunciation by types. In our works, we have used 26 MFCC features as input features.

### C. Dataset

In our works for phonetic classification and detection, we used the dataset named **Japanese Phonetic Database - JPD** [1]. The JPD provides IPA phonetic transcriptions that accurately indicate how Japanese names and words are pronounced in actual speeches, as well as accent codes, for each entry. It includes detail descriptions of 130,000 entries in Japanese.

The advantage of this dataset are: it includes the pitch accent position for each phonetic and it is so meaningful for many speech processing tasks. An example of phonetics collection in JPD is shown below.

Orthographic	Kana	Phonetic	Accent	Pitch Pattern	Remarks
井川	イカワ	[ikawa]	1	HLL(L)	first mora accented
井田	イダ	[ida]	0	LH(H)	accentless
磯貝	イソガイ	[isoŋai]	2	LHLL(L)	second mora accented
鏡	カガミ	[kaŋami]	3	LHH(L)	last mora accented
形	カタチ	[katatçi]	0	LHH(H)	accentless

Fig. 2: JPD phonetics dataset with pitch accent

TABLE I: Japanese phonetic classifier performance of four ANN architectures

	CNN ID	Normal RNN	LSTM	CNN + LSTM
FAR	21.92	17.21	15.93	11.73
FRR	11.29	9.32	9.02	6.23

#### D. Phonetic classification methods

Deep learning is a technique used a lot recently. There are several methods of deep learning that we use there. Recurrent Neural Network (RNN) is a class of artificial neural network where connections between units form a directed graph along a sequence. This allows exhibiting dynamic temporal behavior for a time sequence. We use RNN with 3 layers of LSTM cell. Another deep neural network class we use there is Convolutional Neural Network (CNN), it also find the connectivity between features and help us to improve the accuracy. The architecture of CNN we use there is combination of convolutional layers, maxpool and fully connected layers. Finally we use both RNN and CNN in a combination called deepspeech, they help us increasing the accuracy in this works.

#### E. Result

We tried some models of deep learning architecture mentioned above to training the extended pronunciation recognizer. Following our previous works, two metrics false acceptance rate (**FAR**) and false rejection rate (**FRR**) are used to measure the system performance. The comparison of our Neural Network accuracy with some algorithms is shown in **Table 1**

### III. CONCLUSION AND FUTURE WORKS

#### APPLICATION OF RESEARCH

The research result was applied in **Chatty Pheasant Application** - a service to help non-native learners improving their Japanese pronunciation of **Framgia Inc**<sup>1</sup>

#### ACKNOWLEDGMENT

This research was partially supported by **Framgia Vietnam**. We are thankful to our colleagues who provided expertise that greatly assisted the research, although they may not agree with all of the interpretations provided in this paper.

#### REFERENCES

- [1] Zissman, Marc A., et al. "Automatic dialect identification of extemporaneous conversational, Latin American Spanish speech." Acoustics, Speech, and Signal Processing, 1996. ICASSP-96. Conference Proceedings., 1996 IEEE International Conference on. Vol. 2. IEEE, 1996.
- [2] Harrison, Alissa M., et al. "Implementation of an extended recognition network for mispronunciation detection and diagnosis in computer-assisted pronunciation training." International Workshop on Speech and Language Technology in Education. 2009.
- [3] Ng, Kenney, and Victor W. Zue. "Phonetic recognition for spoken document retrieval." Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on. Vol. 1. IEEE, 1998.
- [4] Zeghidour, Neil, et al. "Learning Filterbanks from Raw Speech for Phone Recognition." arXiv preprint arXiv:1711.01161 (2017).
- [5] Tóth, László. "Phone recognition with hierarchical convolutional deep maxout networks." EURASIP Journal on Audio, Speech, and Music Processing 2015.1 (2015): 25.
- [6] Vaněk, Jan, et al. "A Regularization Post Layer: An Additional Way How to Make Deep Neural Networks Robust." International Conference on Statistical Language and Speech Processing. Springer, Cham, 2017.
- [7] Mohamed, Abdel-rahman, George Dahl, and Geoffrey Hinton. "Deep belief networks for phone recognition." Nips workshop on deep learning for speech recognition and related applications. Vol. 1. No. 9. 2009.
- [8] Tóth, László. "Combining time-and frequency-domain convolution in convolutional neural network-based phone recognition." Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on. IEEE, 2014.
- [9] Graves, Alex, Abdel-rahman Mohamed, and Geoffrey Hinton. "Speech recognition with deep recurrent neural networks." Acoustics, speech and signal processing (icassp), 2013 ieee international conference on. IEEE, 2013.
- [10] Garofolo, John S., et al. "DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1." NASA STI/Recon technical report n 93 (1993).

<sup>1</sup>www.recruit.framgia.vn