ASR

ASEAN

End-to-End

1

This research shows we can integrate linguistic knowledge into the neural network instead of adding more layers or enlarging the model size. The proposed method is universally available for broad tasks for Society 5.0 (such as multilingual speech recognition, disordered speech recognition).

As the most natural way of communication, voice interface with the support of automatic speech recognition (ASR) technology has become crucial in human-computer interaction (HCI) in various devices of today's high-digitized society. Most commercial ASR-enabled products focus on specific popular languages such as English, French, Chinese, Japanese. The speech recognition of less popular languages, such as the ASEAN languages, is still a topic worthy of continued research. Global internationalization raises many real-life situations of multilingual communication, such as regional events, cultural exchanges, festivals.
   The proposed project focused on tackling the problems of the low-resource data and modeling many languages in a single model under the current state-of-the-art End-to-End modeling framework. We also made an in-depth investigation of these problems.

speech recognition  multilingual  articulation  End-to-End

This project will focus on tackling the problems of the low-resource language (e.g., ASEAN languages) and modeling languages as many as we can (hundreds of languages from all language families) in a single model under the current state-of-the-art End-to-End automatic speech recognition (ASR) framework.

Most commercial automatic speech recognition (ASR)-enabled products to focus on specific popular languages such as English, French, Chinese, Japanese. The speech recognition of less popular languages, such as the ASEAN languages, is still a topic worthy of continued research. Global internationalization raises many real-life situations of multilingual communication, such as regional events, cultural exchanges, festivals.

This project will focus on tackling the problems of the low-resource data and modeling many languages in a single model under the current state-of-the-art End-to-End modeling framework. We will also make an in-depth investigation of these problems.

Current work from industry and academia are fascinating with huge self-attentional End-to-End models, and huge multilingual dataset as shown in Figure 1. Recent works from Facebook-AI greatly improved the performance with self-supervised learning and data augmentation from GAN model.
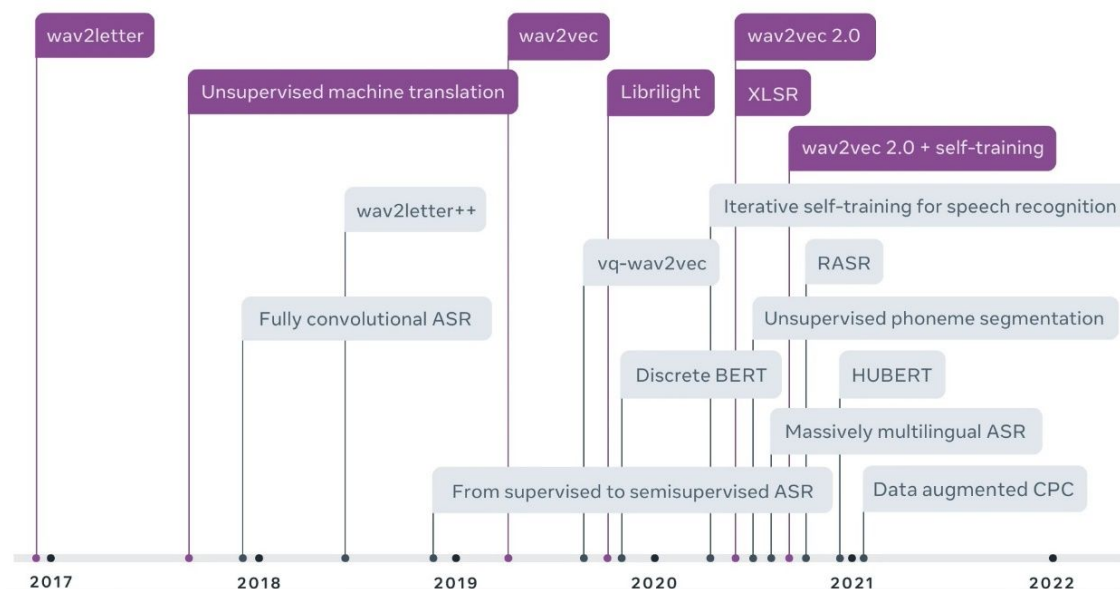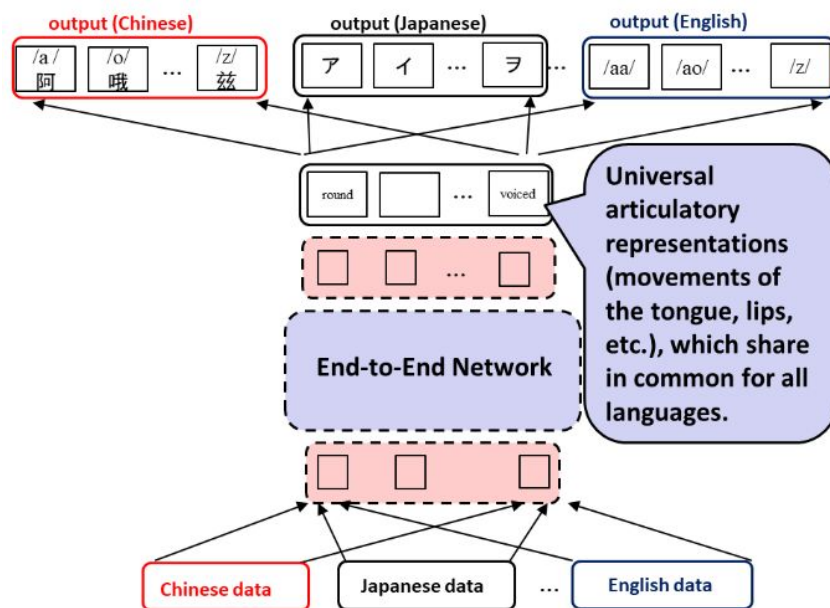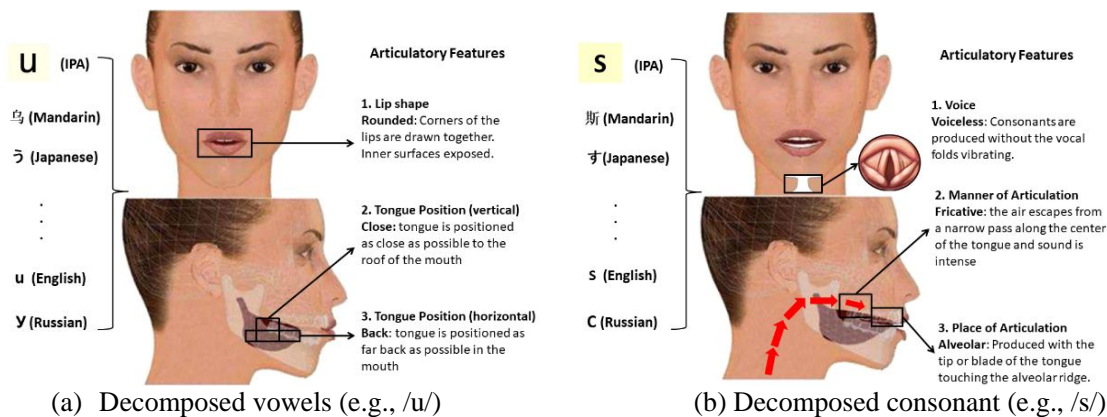


Figure 1: Facebook AI's development timeline of several years of work in speech recognition models, data sets, and training techniques. (cite from https://ai.facebook.com/blog/wav2vec-unsupervised-speech-recognition-without-supervision)

No one can deny the contribution from these works. However, these industry works are quite misleading to common researchers. Someday we can no longer rely on adding more layers or enlarging the model and data size. Our research showed there is another way as alternative: **integrating linguistic knowledge to the neural network.**

In our proposed method, we propose **universal articulatory representation for all of the languages.** We decompose the phones of different languages into "**atom units**": articulatory units as shown in Figure 2 (a, b). Then we can universally model all of the languages together using End-to-End framework. When recognition, we first get articulatory representations, then recover the texts of multilingual languages as shown in Figure 1 (c).



(a) Decomposed vowels (e.g., /u/)

(b) Decomposed consonant (e.g., /s/)

(c) End-to-End multilingual modeling using universal articulatory representation

Figure 2: Universal articulatory representation for all of the languages (a, b) and proposed End-to-End multilingual modeling (c)

This method can effectively represent all of the languages. It has several advantages:

(1) No matter how many languages mix, they can be represented using a universal symbolic set around 20. This technology can make a very compressed neural network output layer for memory limit devices.

(2) Our proposed universal articulatory representation can share knowledge between different languages more effectively. The low-resource languages can also benefit from the other rich sourced languages.

(3) I am so happy to find that the proposed method is universally available for various wide tasks (multilingual mispronunciation detection, multilingual speech recognition, disordered speech recognition, language identification, and speaker diarization). This is our contribution to the community.

In FY2019, I focus on algorithm optimization. I am so happy to find that the proposed method is universally available for various wide tasks (multilingual mispronunciation detection, multilingual speech recognition, disordered speech recognition, language identification, and speaker diarization). Achievements are as follows:

(1) Two international papers were accepted in ICASSP2020, one joint first author and one corresponding author.
(2) One co-authored with an equal contribution in ICME2020.
(3) One first author, and one co-author in Speaker Odyssey2020.
(4) There are also three domestic presentations were reported in ASJ2020.
(5) The patents and book chapters are also included.


In FY2020, I focus on accent speech recognition (English and Chinese), cross-language family speech recognition. Multilingual speech recognition technologies have also been applied to language identification, speaker recognition, disordered speech recognition, and more complex tasks, such as speech translation and adversarial attack. Achievements are as follows:

(1) Multilingual modeling technology has been applied to speaker modeling (1 domestic presentation: IEICE-SP), low-resource transfer learning (1 Interspeech SLIMT2020), and speech translation (NLP2021 presentation), language identification (1 journal paper of IEEE-TASLP), and disordered speech recognition (1 Interspeech2020 with grant honor, 1 O-COCOSDA).

(2) I also find the acoustic modeling unit selection technology can enhance single-language speech recognition with multi-unit (1 invited full paper on 1 Interspeech SLIMT2020, 1 ICASSP2021) and code-switched speech synthesis (1 Interspeech SLIMT2020, 1 ICONIP paper).

(3) Following researches also benefit from the multilingual modeling technologies: speech separation (1 Interspeech2020 with grant honor), adversarial attack (1 IEEE-SLT demo paper), voice-privacy (1 invited report on Interspeech SLIMT2020, 1 Interspeech challenge, 1 ACM-CCS demo), voice activity detection (1 ICASSP2021), Mandarin tone modeling (1 ICASSP2021).

| 1 | 1 | 0 | 0 | |
|---|---|---|---|---|
| P. Shen, X. Lu, S. Li, H. Kawai. | | | | 28 |
| Knowledge Distillation-based Representation Learning for Short-Utterance Spoken Language Identification | | | | 2020 |
| IEEE/ACM Trans. Audio, Speech and Language Process. | | | | 2674 - 2683 |
| DOI 10.1109/TASLP.2020.3023627 | | | | |
| | | | | |

| 24 | 4 | 19 | |
|---|---|---|---|
| S. Li, X. Lu, R. Dabre, P. Shen and H. Kawai | | | |
| Joint Training End-to-End Speech Recognition Systems with Speaker Attributes. | | | |
| ISCA-Odyssey (The Speaker and Language Recognition Workshop) | | | |
| 2020 | | | |

| P. Shen, X. Lu, K. Sugiura, S. Li and H. Kawai. |
|---|
| Compensation on x-vector for short utterance spoken language identification. |
| ISCA-Odyssey (The Speaker and Language Recognition Workshop) |
| 2020 |

| Y. Han, S. Li, Y. Cao, Q. Ma and M. Yoshikawa. |
|---|
| Voice-Indistinguishability: Protecting Voiceprint in Privacy Preserving Speech Data Release. |
| IEEE-ICME |
| 2020 |

Y. Lin, L. Wang, J. Dang, S. Li, and C. Ding,

End-To-End Articulatory Modeling for Dysarthria Articulatory Attribute Detection.

IEEE-ICASSP

2020

H. Shi, L. Wang, M. Ge, S. Li, and J. Dang,

Spectrograms Fusion with Minimum Difference Masks Estimation for Monaural Speech Dereverberation.

IEEE-ICASSP

2020

S. Li, C. Ding, X. Lu, P. Shen and H. Kawai,

End-to-End Articulatory Attribute Modeling for Low-resource Multilingual Speech Recognition,

Acoustical Society of Japan, spring, 2020.

2020

S. Li, X. Lu, R. Dabre, P. Shen and H. Kawai,

Joint Training End-to-End Systems for Speech and Speaker Recognition with Speaker Attributes,

Acoustical Society of Japan, spring, 2020.

2020

P. Shen, X Lu, K Sugiura, S. Li, H Kawai,

Improvement of x-vector for short utterance spoken language identification,

Acoustical Society of Japan, spring, 2020.

2020

N Li, L. Wang, M Unoki, S. Li, R Wang, M Ge, J. Dang,

Robust voice activity detection using a masked auditory encoder based convolutional neural network.

IEEE-ICASSP, 2021

2021

S. Chen, X Hu, S. Li, X Xu,

An investigation of using hybrid modeling units for improving End-to-End speech recognition systems.

IEEE-ICASSP, 2021

2021

H Huang, K Wang, Y. Hu, S. Li,

Encoder-Decoder based pitch tracking and joint model training for Mandarin tone classification.

IEEE-ICASSP, 2021

2021

H. Zhang, S. Li, X. Ma, Y. Zhao, Y. Cao, T. Kawahara,

Phantom in the Opera: Effective Adversarial Music Attack on Keyword Spotting Systems.

IEEE-SLT, 2021

2021

K. Soky, S. Li, M. Mimura, C. Chu, T. Kawahara,

Comparison of End-to-End Models for Joint Speaker and Speech Recognition

IEICE-SP, 2021.

2021

K. Soky, S. Li, T. Kawahara, S. Seng,

Multilingual transformer training for Khmer automatic speech recognition

Interspeech 2020 Satellite Workshop (SLTMS2020)

2020

S. Shimizu, C. Chu, S. Li, S. Kurohashi,

End-to-End Speech Translation with Cross-lingual Transfer Learning

NLP, 2021.

2020

S. Guo, L. Wang, S. Li, J. Zhang, C. Gong, Y. Wang, J. Dang, K. Honda

Effectively Synthesizing Code-switched Speech Using Highly Imbalanced Mix-lingual Data and mask embedding

Interspeech 2020 Satellite Workshop (SLTMS2020)

2020

H. Zhang, S. Ueno, M. Mimura, S. Li, W. Zhang, T. Kawahara,

A Mixture of Character and Word End-to-End System for Keyword Spotting

Interspeech 2020 Satellite Workshop (SLTMS2020)(full paper).

2020

S. Guo, L. Wang, S. Li, J. Zhang, C. Gong, Y. Wang, J. Dang, K. Honda.

Effectively Synthesizing Code-switched Speech Using Highly Imbalanced Mix-lingual Data

In Proc. ICONIP, 2020

2020

Y. Lin, L. Wang, S. Li, J. Dang, and C. Ding,

Staged Knowledge Distillation for End-to-End Dysarthric Speech Recognition and Speech Attribute Transcription

In Proc. INTERSPEECH, 2020 (Travel Granted by ISCA).

2020

A. Thida, N. Han, S. Oo, S. Li and C. Ding,

VOIS: The First Speech Therapy App in the World for Myanmar Hearing-Impaired Children.

In Proc. O-COCOSDA, 2020.

2020

Y. Han, Y. Cao, S. Li, Q. Ma, M. Yoshikawa.

Voice-Indistinguishability: Protecting Voiceprint in Privacy-Preserving Speech Data Release,

Interspeech 2020 Satellite Workshop (SLIMS2020) (invited report).

2020

Y. Han, Y. Cao, S. Li, Q. Ma, M. Yoshikawa.

Voice-Indistinguishability: Protecting Voiceprint with Differential Privacy under an Untrusted Server.

ACM conference on Computer and Communications Security (CCS), demo, 2020.

2020

Y. Han, S. Li, Y. Cao, M. Yoshikawa,

System Description for Voice Privacy Challenge (Kyoto Team).

In special session of INTERSPEECH 2020 (VoicePrivacy challenge 2020).

2020

| | |
|---|---|
| H Shi, L. Wang, S. Li, C. Ding, M Ge, N Li, J. Dang, and H Seki. | |
| Singing Voice Extraction with Attention based Spectrograms Fusion. | |
| In Proc. INTERSPEECH, 2020 (Travel Granted by ISCA). | |
| 2020 | |

1

| | |
|---|---|
| X. Lu, S. Li, M Fujimoto | 2020 |
| Springer Singapore | 18 |
| Automatic speech recognition | |

4

| | | |
|---|---|---|
| | | |
| 2019-163555 | 2019 | |

| | | |
|---|---|---|
| | | |
| 2020-059962 | 2020 | |

| | | |
|---|---|---|
| | , , | |
| 2019-086005 | 2019 | |

| | | |
|---|---|---|
| | | |
| 2019-051008 | 2019 | |

0

O

| | Tianjin University | Xinjiang University | H think RoyalFlush AI | |
|---|---|---|---|---|