

UNIT SELECTION SYNTHESIS BASED DATA AUGMENTATION FOR FIXED PHRASE SPEAKER VERIFICATION

Houjun Huang^{1,2}, Xu Xiang¹, Fei Zhao¹, Shuai Wang², ✉ Yanmin Qian²

¹ AISpeech Ltd, Suzhou China

²MoE Key Lab of Artificial Intelligence, AI Institute SpeechLab, Department of Computer Science and Engineering Shanghai Jiao Tong University, Shanghai, China
 {houjun.huang, xu.xiang, fei.zhao01}@aispeech.com, {feixiang121976, yanminqian}@sjtu.edu.cn

ABSTRACT

Data augmentation is commonly used to help build a robust speaker verification system, especially in limited-resource case. However, conventional data augmentation methods usually focus on the diversity of acoustic environment, leaving the lexicon variation neglected. For text dependent speaker verification tasks, it's well-known that preparing training data with the target transcript is the most effectual approach to build a well-performing system, however collecting such data is time-consuming and expensive. In this work, we propose a unit selection synthesis based data augmentation method to leverage the abundant text-independent data resources. In this approach text-independent speeches of each speaker are firstly broke up to speech segments each contains one phone unit. Then segments that contain phonetics in the target transcript are selected to produce a speech with the target transcript by concatenating them in turn. Experiments are carried out on the AISHELL Speaker Verification Challenge 2019 database, the results and analysis shows that our proposed method can boost the system performance significantly.

Index Terms— speaker verification, data augmentation, unit selection synthesis, x-vector

1. INTRODUCTION

Speaker verification (SV) aims to confirm the claimed identity given his/her speech. Considering the constraint on the speech content, the SV task can be further categorized into two classes: text-dependent and text-independent. The former requires the same content for the enrollment and test utterances, while the latter doesn't.

Traditional speaker recognition systems are based on statistical models such as Gaussian Mixture Model-Universal Background Model (GMM-UBM) [1]. The i-vector [2] system projects the GMM super-vector to a lower-dimensional

and more speaker-discriminative vector. Recently, utterance-level deep speaker embedding methods such as x-vector [3, 4], have shown better performance than i-vector on many standard speaker recognition data-sets.

Although the SV performance has been improved greatly over the recent years, further upgrading its robustness in real applications is still a challenge task due to the complex environment. Data augmentation is conventionally adopted in building a SV system to improve its robustness. Snyder *et al.* in [3, 4] manually employed additive noises and reverberation to original speech segments in the training set to train a robust embedding extractor and then extract "clean" and "noisy" embeddings to train probabilistic linear discriminant analysis (PLDA) for both i-vector and x-vector. SpecAugment [5, 6] also shows promising results on the speaker verification task. On the other hand, researchers also apply deep generative models such as generative adversarial network (GAN) and variational autoencoder (VAE) to generate x-vector embeddings directly to train a robust PLDA [7, 8].

Despite the effectiveness exhibited by the data augmentation methods above, they only consider the variation of acoustic environment. Text variation should also be considered to build a robust SV system, especially for text-dependent tasks. In practice, to achieve the state-of-the-art performance, text-dependent SV requires the same set of text to be spoken during the training stage and the test stage. [9, 10, 11]. When there is no or limited training data with the designated phrase to build a text-dependent SV system, those data augmentation approaches couldn't help to get promising performance. To improve the system performance from this aspect, in this work, we propose a novel method which generates new speech containing designated phrase from text-independent database using the *unit selection synthesis* [12, 13, 14].

2. UNIT SELECTION SYNTHESIS BASED DATA AUGMENTATION

In the case that limited or no text-dependent training data is available to build a text-dependent SV system, generating

Yanmin Qian is the corresponding author. This work was supported by the National Key R&D Program of China (No. 2018YFB1004602) and the China NSFC project (No. 62071288).

speech utterances with the fixed phases of more speakers is the most effectual solution. If a text-independent database with a large number of units is available, speaker discriminative synthesized speeches can be produced by concatenating the wave-forms of units selected from speeches of each speaker [12, 13, 14].

Here, we use the Chinese wake-up word "ni hao mi ya" as an example to describe how to carry out the proposed approach. In this work, we treat each Chinese character as the phonetic unit. Thus, "ni hao mi ya" can be converted to a phonetic unit sequence "ni"- "hao"- "mi"- "ya." The unit selection synthesis based augmentation is shown in Figure 1, which contains the following steps,

1. For each speaker in the text-independent database, we chunk his speech into segments which only contain single characters.
2. Select all the segments which contain desired phonetic units to generate phonetic unit libraries.
3. For each speaker, synthesize new speech containing the target transcript by concatenating the sampled units from the phonetic unit libraries.

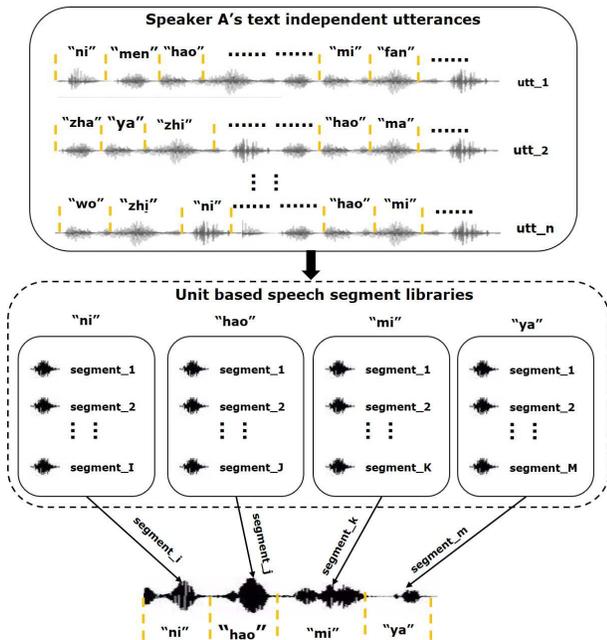


Fig. 1. Generate "ni hao, mi ya" using unit selection synthesis for a speaker in the text-independent database

The standard unit selection synthesis system aims to produce more natural-sounding synthesized speeches, while the proposed approach aims at generating utterances with fix phonetic content that are speaker discriminative and variable to train more robust text-dependent SV systems. Speech segments are randomly selected from the same speaker for each

unit and no flexible join technique is used in concatenating to ensure diversity of synthesized data. In our experiments, N_i utterances will be synthesized for speaker i where N_i is the max size of his/her unit based speech segment libraries.

3. EXPERIMENTS

3.1. Experimental data

In this work, we use the AISHELL2¹ [15] as text-independent training data, AISHELL-wakeup² [11] as text-dependent training data, and AISHELL-2019B-eval dataset³ [11] as text-dependent test set.

The AISHELL2 is a 1000-hour Mandarin Chinese Speech Corpus that contains 1003664 close-talk utterances from 1991 speakers. The speech utterance contains several domains, including keywords, voice command, smart home, autonomous driving, industrial production, etc. The recording was put in a quiet room environment using an iOS system mobile phone.

The AISHELL-wakeup database has 3,936,003 utterances from 254 speakers. The content of utterances covers two wake-up words: "ni hao, mi ya" in Chinese and "Hi, Mia" in English. The recording process was put in a real smart home environment where one close-talking microphone was placed 25cm away from the speaker, six 16-channel circular microphone arrays were placed around the person with a distance including 1m, 3m and 5m from the speaker and the noise source was randomly placed close to one of the microphone arrays. The 993083 mono channel wave-files of the Chinese wake-up word "ni hao, mi ya" are chosen to train the text-dependent model in our experiments.

The AISHELL-2019B-eval contains recordings of 86 speakers with Chinese wake-up word "ni hao, mi ya". The room setting and recording devices are the same as that of AISHELL-wakeup. Utterances of the last 44 people are selected as the test set since they are more challenging [11]. This corpus has two tasks: close-talking enrollment task (utterances from the close-talking mic are used for enrollment) and far-field enrollment task (utterances from one 16-channel circular microphone array which is 1m away from the speaker are used for enrollment). The testing data for both tasks are far-field utterances recorded with 16-channel circular microphone arrays.

3.2. Experimental setups

The frame-level alignments of the above databases are generated using the official Kaldi [16] AISHELL2 speech recognition recipe(s5) [15], and voice activity detection (VAD) labels of them are generated based on their alignments.

¹ AISHELL2 is publicly available at http://www.aishelltech.com/aishell_2.

² AISHELL-wakeup is available at <http://openslr.org/85/>.

³ Speech data of AISHELL-2019B-eval and trial files are available at <http://openslr.org/85/>.

347706 utterances of 1986 speakers are generated from the AISHELL2 database using the proposed approach, and we call this data-set as "AISHELL2-aug" in the following sections. After the data augmentation method proposed in [3, 4] is adopted, the number of speech utterances of AISHELL2, AISHELL-wakeup and AISHELL2-aug are extended to 4013020, 3972332 and 3824766. 40-dimensional fbank features are extracted from these databases with a frame shift of 10ms and a window width of 25ms. VAD is employed to filter out non-speech frames. Mean normalization is then applied over a sliding window of up to 300 frames. Finally, SpecAugment [5] is applied to the fbank features.

The x-vector system [3, 4] is used in our experiments. The architecture of the speaker-discriminative TDNN is illustrated in Table 1. The T in the stats-pool layer corresponds to the frame number of the input features. X-vector embeddings are extracted at layer segment6, before the projection layer. The N in the projection layer corresponds to the number of training speakers. Additive angular margin loss [17, 18] with $m=0.2$ and $s=32.0$ is used as the projection layer since it has shown better performance than other losses. For the enrollment or test utterances that have 16 channels from a 6-channel circular microphone array, we adopt the strategy of speaker embedding level averaging as do in [11].

Table 1. Architecture of the x-vector

Layer	Layer context	Total context	Input*output
frame1	[t-2, t+2]	5	200*256
frame2	t-2, t, t+2	9	768*256
frame3	t-3, t, t+3	15	768*256
frame4	t	15	256*256
frame5	t	15	256*512
stats-pool	[0,T)	T	512T*1024
segment6	0	T	1024*256
projection	0	T	256*N

AISHELL2 is firstly used to train the text-independent x-vector model. As is shown in Figure 2, when we train a text-dependent model, the parameters of layers before the projection layer are initialized by the text-independent x-vector model. To test the effectiveness of the proposed approach, three text-dependent x-vector models are trained with AISHELL-wakeup, AISHELL2-aug and AISHELL-wakeup + AISHELL2-aug, respectively.

Cosine similarity serves as back-end scoring method during testing. We report results in Equal error rate (EER), minimum detection cost for $P(\text{tar}) = 0.01$ (minDCF).

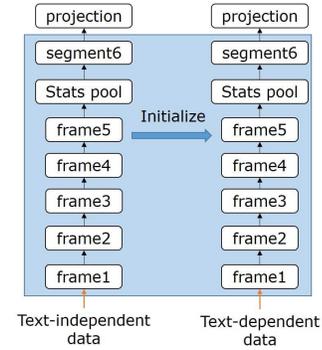


Fig. 2. Initialize text-dependent x-vector model with text-independent model

3.3. Experimental results

Experiment results on close-talking and far-field enrollment task are shown in Table 2 and Table 3 respectively.

Table 2. Performance on close-talking enrollment task of x-vector models trained with different data. *AISHELL2-aug* is the training data produced by the proposed approach.

train data set	EER(%)	minDCF
AISHELL2	6.978	0.616
AISHELL-wakeup	5.796	0.606
AISHELL-wakeup+AISHELL2	4.386	0.446
AISHELL2-aug	4.087	0.405
AISHELL-wakeup+AISHELL2-aug	3.019	0.282

Table 3. Performance on far-field enrollment task of x-vector models trained with different data. *AISHELL2-aug* is the training data produced by the proposed approach.

train data set	EER(%)	minDCF
AISHELL2	5.993	0.466
AISHELL-wakeup	4.598	0.409
AISHELL-wakeup+AISHELL2	3.274	0.335
AISHELL2-aug	3.673	0.319
AISHELL-wakeup+AISHELL2-aug	2.562	0.224

Compared to the text-independent x-vector model, when there is no text-dependent resource to build the SV system, AISHELL2-aug generated by the proposed method trained text-dependent x-vector model achieves a relative performance improvement of 41.4% in EER and 38.7% in EER on close-talking enrollment task and far-field enrollment task respectively. Compared to the AISHELL-wakeup+AISHELL2

trained text-dependent x-vector model, when there is limited text-dependent resource to build the SV system, AISHELL-wakeup+AISHELL2-aug trained model obtains further relative performance improvement of 31.16% in EER and 21.74% in EER on close-talking enrollment task and far-field enrollment task, respectively.

3.4. Analysis

The mismatch of speech content between training and test data will induce a severe degradation of SV performance. The proposed method aims to produce training speech utterances whose contents match the text-dependent test set. This approach is effective only if the synthesized speeches contain a phonetic context of "ni"-“hao”-“mi”-“ya” and be speaker discriminative.

A deep-neural network(DNN) based keyword spotting system is firstly tested on AISHELL-wakeup, AISHELL2-aug and AISHELL2 (speeches in AISHELL-wakeup or AISHELL2-aug are positive samples, speeches in AISHELL2 are negative samples). The DNN with 7 hidden layers and 256 nodes per hidden layer is pre-trained with about 5000h speeches. The DNN has 411 output labels: 409 Chinese characters, a silence label and a Filler label for music and noise. 40-dimensional fbank features are extracted with a frame shift of 20ms and a window width of 30ms and then 5 future frames and 5 frames in the past are stacked to predict posterior probabilities for each output label using the DNN model. The posterior handling module proposed in [19] combines the label posteriors produced every frame into a confidence score used for detection. Figure 3 shows the performance when speeches in AISHELL-wakeup or AISHELL2-aug are chosen as positive samples. Results are demonstrated in the form of Receiver Operating Characteristic (ROC) curves, where the false reject rate(that is, a key phrase is present but a negative decision is given, FRR) is on the Y-axis and the false alarm rate (that is, a key-phrase is not present, but a positive decision is made, FAR) is on X-axis. The ROC is obtained by sweeping through confidence thresholds.

As is shown in Figure 3, when FAR is larger than 0.01, speeches in AISHELL2-aug could achieve very similar FRR with speeches in AISHELL-wakeup. This means the synthesized speech does contain a phonetic context of "ni"-“hao”-“mi”-“ya”. As the synthesised speeches are not natural-sounding, when FAR drops to 0.001, FRR of AISHELL-aug is 0.164 while that of AISHELL-wakeup is 0.087. In our future work, we will try to produce or select more natural-sounding synthesized speeches to train the x-vector model.

The text-independent x-vector model trained with the Voxceleb2 database in our previous work [18] is used to extract embeddings from audios of AISHELL2-aug. 50 speakers are randomly chosen, and their embeddings are shown in Figure 4 using t-SNE. Figure 4 shows the good speaker discriminative property of synthesized speeches.

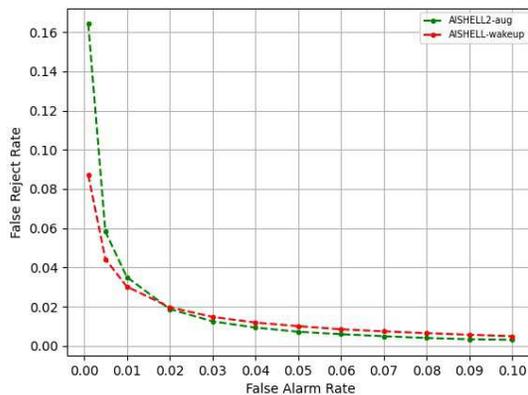


Fig. 3. ROC curves when speeches in AISHELL-wakeup or AISHELL2-aug are chosen as positive samples

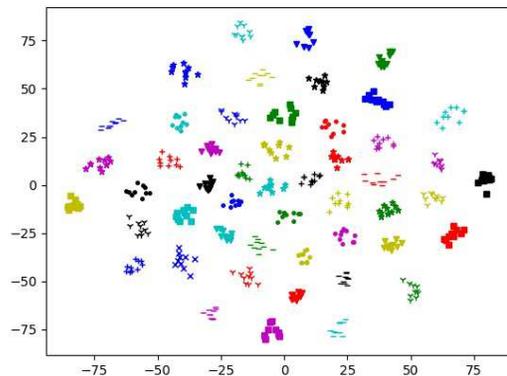


Fig. 4. t-SNE visualization of embeddings from 50 random speakers of AISHELL2-aug, samples with the same shape and color are from the same speaker

4. CONCLUSIONS AND FUTURE WORK

In this work, we proposed a unit selection synthesis based data augmentation for text-dependent speaker verification. The proposed method enables us to leverage the rich text-independent data to fast generate new speech with the desired transcript, which leads to a better-performing and more robust text-dependent speaker verification system. This strategy can reduce the development period and cost dramatically. Experiments on the AISHELL-2019B-eval corpus shows that the proposed approach could achieve a relative performance improvement of about 40% in both EER and minDCF.

In the future work, we will focus on synthesizing more natural-sounding and variable speech to further increase the robustness of speaker verification systems.

5. REFERENCES

- [1] Douglas A Reynolds, Thomas F Quatieri, and Robert B Dunn, "Speaker verification using adapted gaussian mixture models," *Digital signal processing*, vol. 10, no. 1-3, pp. 19–41, 2000.
- [2] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.
- [3] David Snyder, Daniel Garcia-Romero, Daniel Povey, and Sanjeev Khudanpur, "Deep neural network embeddings for text-independent speaker verification.," in *Interspeech*, 2017, pp. 999–1003.
- [4] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5329–5333.
- [5] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," *arXiv preprint arXiv:1904.08779*, 2019.
- [6] Shuai Wang, Johan Rohdin, Oldřich Plchot, Lukáš Burget, Kai Yu, and Jan Černocký, "Investigation of specAugment for deep speaker embedding learning," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7139–7143.
- [7] Yexin Yang, Shuai Wang, Man Sun, Yanmin Qian, and Kai Yu, "Generative adversarial networks based x-vector augmentation for robust probabilistic linear discriminant analysis in speaker verification," in *2018 11th International Symposium on Chinese Spoken Language Processing (ISCSLP)*. IEEE, 2018, pp. 205–209.
- [8] Zhanghao Wu, Shuai Wang, Yanmin Qian, and Kai Yu, "Data augmentation using variational autoencoder for embedding based speaker verification.," in *INTER-SPEECH*, 2019, pp. 1163–1167.
- [9] Yexin Yang, Shuai Wang, Xun Gong, Yanmin Qian, and Kai Yu, "Text adaptation for speaker verification with speaker-text factorized embeddings," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6454–6458.
- [10] Xiaoyi Qin, Danwei Cai, and Ming Li, "Far-field end-to-end text-dependent speaker verification based on mixed training data with transfer learning and enrollment data augmentation.," in *INTER-SPEECH*, 2019, pp. 4045–4049.
- [11] Xiaoyi Qin, Hui Bu, and Ming Li, "Hi-mia: A far-field text-dependent speaker verification database and the baselines," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7609–7613.
- [12] Andrew J Hunt and Alan W Black, "Unit selection in a concatenative speech synthesis system using a large speech database," in *1996 IEEE International Conference on Acoustics, Speech, and Signal Processing Conference Proceedings*. IEEE, 1996, vol. 1, pp. 373–376.
- [13] Alan W Black and Paul A Taylor, "Automatically clustering similar units for unit selection in speech synthesis.," 1997.
- [14] Alistair Conkie, "Robust unit selection system for speech synthesis," in *137th meeting of the Acoustical Society of America*, 1999, p. 978.
- [15] Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu, "Aishell-2: transforming mandarin asr research into industrial scale," *arXiv preprint arXiv:1808.10583*, 2018.
- [16] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hanemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, et al., "The kaldi speech recognition toolkit," in *IEEE 2011 workshop on automatic speech recognition and understanding*. IEEE Signal Processing Society, 2011, number CONF.
- [17] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [18] Xu Xiang, Shuai Wang, Houjun Huang, Yanmin Qian, and Kai Yu, "Margin matters: Towards more discriminative deep neural network embeddings for speaker recognition," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2019, pp. 1652–1656.
- [19] Guoguo Chen, Carolina Parada, and Georg Heigold, "Small-footprint keyword spotting using deep neural networks," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014.