

TEACH AN ALL-ROUNDER WITH EXPERTS IN DIFFERENT DOMAINS

Zhao You¹, Dan Su¹, Dong Yu²

¹Tencent AI Lab, Shenzhen, China

²Tencent AI Lab, Bellevue, WA, USA

{dennisyou, dansu, dyu}@tencent.com

ABSTRACT

In many automatic speech recognition (ASR) tasks, an ideal model has to be applicable over multiple domains. In this paper, we propose to teach an all-rounder with experts in different domains. Concretely, we build a multi-domain acoustic model by applying the teacher-student training framework. First, for each domain, a teacher model (domain-dependent model) is trained by fine-tuning a multi-condition model with domain-specific subset. Then all these teacher models are used to teach one single student model simultaneously. We perform experiments on two predefined domain setups. One is domains with different speaking styles, the other is near-field, far-field and far-field with noise. Moreover, two types of models are examined: deep feedforward sequential memory network (DFSMN) and long short term memory (LSTM). Experimental results show that the model trained with this framework outperforms not only multi-condition model but also domain-dependent model. Specially, our training method provides up to 10.4% relative character error rate improvement over baseline model (multi-condition model).

Index Terms— multi-domain, all-rounder, speech recognition, teacher-student training, knowledge distillation

1. INTRODUCTION

Thanks to deep learning approaches [1, 2], great progress has been made in automatic speech recognition performance. Although deep neural networks have superior robustness over GMM systems on different conditions such as speaker, recording channel and acoustic environment [3], domain robustness is still a challenging problem. First, it is impractical to train a single model with good performance across all domains. Second, when encountering a new domain task, the models trained with other domain data are usually difficult to transfer knowledge to this new task.

Much effort has been devoted to solve these problems. One of the most effective and straightforward approach is multi-condition training or multi-style training. Multi-condition training can trace back to [4], and has been shown to reduce mismatch for different noise conditions [5]. Deep neural networks work particularly well with multi-condition

training due to the large model capacity [6, 7, 8]. Empirically, when sufficient amount of target domain data is available, performing fine-tuning on the multi-condition model with target domain data can produce excellent performance for the target domain specifically.

Recently, domain adaption based methods have been proposed for domain robustness. Domain adaptation refers to the task of adapting models trained on one domain or mixed domains to the target domain. Previous study mainly focuses on the scenario where with limited target domain training data [9], one has to make a trade-off between test accuracy and the number of adaptation parameters. Further, domain adaptation for ASR is particularly difficult considering the mismatch in speaking styles, noise types, and room acoustics etc [10].

Methods which augment DNN input with i-vector feature or speaker code have been developed against the speaker variations and channel variations. However, it is difficult for these methods to deal with other domain changes such as speaking style variations. For example, in spontaneous speech, the speaking rate is highly inconsistent and the articulation is highly variable, which are typically not observed in read speech. Moreover, the above approach has drawbacks which make it unsuitable for transferring knowledge across multiple domains[11].

In this paper, we propose a multi-domain teacher-student training framework for teaching an all-rounder with experts in different domains. First, we train teacher models for each domain by fine-tuning a multi-condition model with the domain specific subset. Then, we exploit all these teacher models to train one single student model simultaneously. Consequently, the student model is an all-rounder across different domains. Our experimental results show that given sufficient amount of data for several domains the proposed method provides up to 10.4% relative character error rate improvement over baseline models.

The rest of the article is organized as follows. The multi-domain teacher-student training method is described in Section 2. The experimental configuration is described in Section 3. We report the experimental results in Section 4 and conclude the paper in Section 5.

2. MULTI-DOMAIN TEACHER-STUDENT TRAINING

2.1. Teacher-student training

In the teacher-student training framework, [12, 13, 14] have shown that it is possible to train an student model to match the output distribution of a teacher model. Specially, the student model can be learned via single teacher network or multiple teacher networks. Details of the two learning methods will be discussed in the following section.

2.1.1. single teacher network

Teacher-student training is a general method for compressing acoustic model by minimizing the Kullback-Leibler divergence (KLD) between the output posterior distributions of the teacher model and the student model. That is, minimizing the loss function $L_{KLD}(\theta)$ defined as

$$L_{KLD}(\theta) = - \sum_l p_t(l|x) \log p_s(l|x) \quad (1)$$

where $p_t(l|x)$ is the posterior probability of label l given the input feature x computed by the teacher model, and $p_s(l|x)$ is that computed by the student model. θ denotes the parameters of the student model.

However, the valuable knowledge transferred from one teacher network to train the student network is limited. Thus, developing training approaches which use multiple teachers is needed.

2.1.2. multiple teacher networks

In the case of multiple teacher networks, the performance of the student network is improved by leveraging information from multiple teachers. In [15, 16], the outputs of multiple teacher networks are combined by weighted ensembles of posteriors from each teacher network as

$$p_t(l|x) = \sum_{k=1}^N w_k p_{tk}(l|x) \quad (2)$$

where N is the number of teacher networks. $p_{tk}(l|x)$ is the posterior of the k -th teacher network. w_k is the interpolation weight.

Though this approach allows the students to learn ensemble distribution created by multiple teachers, the characteristic of each teacher network is weakened by the interpolation process. In that case, the student network can not obtain the most professional characteristic of each domain from the teacher network's knowledge.

2.2. Multi-domain training

To overcome the shortcomings of methods above, we propose a multi-domain teacher-student learning algorithm. In this

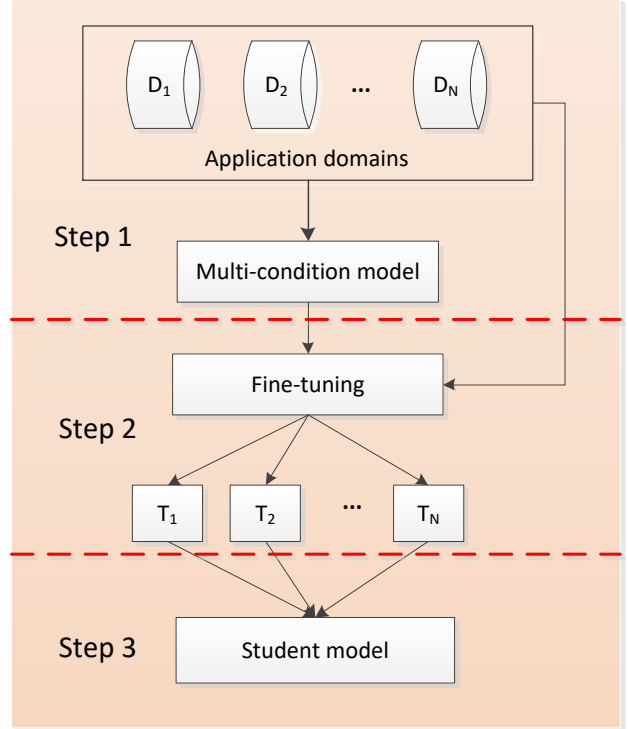


Fig. 1. Training steps of multi-domain teacher-student learning.

method, each teacher network is trained with a domain specific data, which can be viewed as an expert of this domain and transfers the most professional characteristic of each domain. The multi-domain teacher-student training framework is shown in figure 1. The process of training one student network from multiple teacher networks includes three steps:

1. We pool data from multiple application domains. D_n denotes the n -th domain. Then, we train a multi-condition model with minibatches samples which are chosen randomly from the pooled set.

2. Domain-dependent teacher models are produced by fine-tuning the multi-condition model with domain-dependent data respectively. T_n denotes the n -th teacher model which is trained with the n -th domain data.

3. The proposed model is learned from these N domain-dependent teacher models. During the training process, samples in one minibatch are chosen randomly from the mixed data set, and may come from different domains. The train process exploits each sample for training by using the soft targets produced from its corresponding domain-dependent teacher model as in equation (3). Let $p_t^d(l|x)$ be the soft targets produced by domain specific data. $\delta_t(l)$ denotes the hard labels and w_{hard} is its weight. Thus, $p_t(l|x)$ can be viewed as a linear interpolation of hard labels and soft labels.

$$p_t(l|x) = (1 - w_{hard})p_t^d(l|x) + w_{hard}\delta_t(l) \quad (3)$$

Table 1. Total length (approx.hours) of each application domain for mixed speaking style corpus.

speaking-style	train	dev	test
Read	2k	1	1
Lect	2k	1	1
Spon	2k	1	1

Table 2. Total length (approx.hours) of each application domain for mixed environment corpus.

environment	train	dev	test
Near	2k	1	1
Far	2k	1	1
FarNoise	2k	1	1

3. EXPERIMENTAL SETUP

3.1. Training setup

The feature vectors used in all the experiments are 40-dimensional log-mel filterbank energy features appended with the first and second order derivatives. Log-mel filterbank energy features are computed with a 25ms window and shifted every 10ms. We stack 8 consecutive frames and sub-sample the input frames with 3. A global mean and variance normalization is applied for each frame. All the experiments are based on the CTC learning framework. We use the CI-syllable-based acoustic modeling method [17] for the CTC learning. The target labels of CTC learning are defined to include 1394 Mandarin syllables, 39 English phones and a blank. Character error rate (CER) results are measured on the test sets. Decoding is performed with a beam search algorithm by using the weighted finite-state transducers (WFSTs).

3.2. Datasets

Our training corpus consists of a variety of application domains, all in Mandarin. In this work, we evaluate our approach on two kinds of domain setups. The number of utterances and the total length of each application domains are shown in Table 1 and Table 2.

The first setup focuses on different speaking styles. We experiment on a mixed dataset with 3 kinds of different speaking styles, including read speech, lecture speech and spontaneous speech. We refer them as Read, Lect and Spon respectively. Each speaking style contains 2000 hours speech. The second setup focuses on the variation of environment, including near-field speech, simulated far-field speech and simulated far-field noisy speech. We refer them as Near, Far and FarNoise respectively. Together they contain a total of 6000

Table 3. CER (%) of 3 different training methods with DFSMN models. Results are with corpus of 3 different speaking styles.

DFSMN	test-Read	test-Spon	test-Lect
Baseline	17.37	23.18	15.92
T1 (Read)	16.73	34.77	26.77
T2 (Spon)	26.95	21.98	21.57
T3 (Lect)	29.03	29.54	15.88
student model	16.37 (-5.8%)	20.76 (-10.4%)	15.13 (-5.0%)

hours speech. The far-field speech is generated using the image method described in [18]. A set of simulated room impulse responses (RIRs) are created with different rectangular room sizes, speaker positions and microphone positions, as proposed in [19]. Each environment contains 2000 hours speech. We hold out about 0.5% (1 hours) as a development set for frame accuracy evaluation. Each test set includes 1k utterances which is about 1 hours of audio.

3.3. Acoustic Model

We present our work with DFSMN and LSTM acoustic models. The LSTM system uses 7 LSTM layers of 1024 cells, each with a recurrent projection layer of 512 units. For DFSMN model, we use 30 DFSMN components [20]. The look-back order and lookahead order of each memory block is 5 and 1 respectively, and the strides are 2 and 1 respectively. For stable CTC learning, we clip gradients to [-1.0, 1.0]. We use the Kaldi [21] toolkit to train models and all models are trained in a distributed manner using BMUF [22] optimization with 8 Tesla P40 GPUs.

4. EXPERIMENTAL RESULTS

In this work, we evaluate the performance of the proposed method on several large vocabulary Mandarin speech recognition tasks including near-field speech and far-field speech as described in section 3.2.

4.1. Mixed speaking style corpus

For the first set of experiments, we validate the effectiveness of the proposed method by dealing with mixed speaking style speech. To find an appropriate value for w_{hard} in equation (3), we randomly select about 25 % data from each domain to constitute a training set and perform the experiment. We find that $w_{hard} = 0.8$ achieves best performance on the development set. Thus, we set $w_{hard} = 0.8$ for all the experiments.

Table 3 shows the performance comparison on mixed speaking style corpus with DFSMN acoustic models. Line 2

Table 4. CER (%) of 3 different training methods with LSTM models. Results are with corpus of 3 different speaking styles.

LSTM	test-Read	test-Spon	test-Lect
Baseline	17.49	21.09	15.49
T1 (Read)	17.23	29.78	22.11
T2 (Spon)	23.63	19.68	19.75
T3 (Lect)	25.25	25.92	14.58
student model	16.79 (-4.0%)	19.85 (-5.9%)	14.99 (-3.2%)

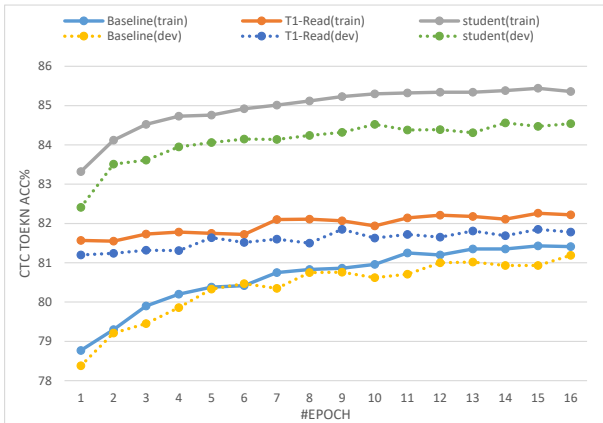


Fig. 2. Frame accuracy on *Read* train and development set against the number of epochs through the training dataset for 3 different training methods.

presents the results of the baseline system (multi-condition model). The following three lines present results of 3 teacher models trained on *Read*, *Lect* and *Spon* data respectively. The last line presents results of the student model which is learned by using 3 teacher models. Compared with the baseline system, the teacher models work well when the test domains match the model domains while poorly when the domains mismatch. Specially, the results clearly show that the student model performs best compared with other domain-adaptation (teacher) models. Finally, our proposed method achieves up to 10.4% relative CER improvement over the baseline model on the test-Spon test set. Table 4 shows the corresponding results of LSTM networks. It can be concluded that the student network still outperforms the baseline system on all the domain tests. However, the gain is smaller compared with DFSMN models.

Figure 2 shows the frame accuracy of different models. As shown, the "student" model produces the highest frame accuracy compared with both the baseline model and the "teacher" model. Since the "student" model is trained on the linear interpolation of the hard targets and soft targets, this suggests that the soft targets role as a significant contribution to train the student model of high accuracy.

Table 5. CER (%) of 3 different training methods with DFSMN models. Results are with corpus of 3 different far-field environments.

DFSMN	test-Near	test-Far	test-FarNoise
Baseline	17.08	26.04	46.28
student model	15.76 (-7.7%)	23.44 (-10.0%)	42.47 (-8.2%)

Table 6. CER (%) of 3 different training methods with LSTM models. Results are with corpus of 3 different far-field environments.

LSTM	test-Near	test-Far	test-FarNoise
Baseline	17.04	26.36	44.03
student model	16.74 (-1.8%)	25.60 (-2.9%)	43.00 (-2.3%)

4.2. Mixed near-field and far-field corpus

To comprehensively validate the effectiveness of our proposed method, we also investigate the performance on the mixed near-field and far-field corpus. The multi-condition model, student model and teacher models have the same configuration with that used in Part 4.1. Table 5 shows that the DFSMN student model trained by our proposed method also significantly outperforms the baseline system. In particular, our training method provides up to 10% relative CER improvement over the baseline model on the test-Far test set. This shows that our proposed method can achieve consistent improvements, no matter on a mixed speaking styles corpus or a mixed near-field and far-field corpus. Table 6 shows the corresponding results of LSTM student model. An observation is that the improvement on LSTM models is smaller compared with DFSMN models.

5. CONCLUSIONS AND FUTURE WORKS

In this paper, we propose a multi-domain teacher-student training method for teaching an all-rounder with experts in different domains. We explore this method for acoustic modeling on two different tasks. We find that the model trained by this method not only outperforms multi-condition model but also outperforms the domain-dependent model produced by fine-tuning a multi-condition model with the target domain data set. Table 3 shows that our method achieves up to 10.4% relative CER improvement over the baseline model.

[16] has shown that combining intermediate representations of multiple teacher networks can significantly improve the student network's performance. Thus, we will explore this training strategy to improve the performance of LSTM models in the future work.

6. REFERENCES

- [1] G. E. Dahl, D. Yu, L. Deng, and A. Acero, "Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," in *IEEE Transactions on audio, speech, and language processing*. IEEE, 2012, vol. 20, p. 3042.
- [2] D. Yu and J. Li, "Recent progresses in deep learning based acoustic models," in *IEEE/CAA Journal of Automatica Sinica*. IEEE, 2017, vol. 4, p. 396409.
- [3] Y. Huang, D. Yu, C. Liu, and Y. Gong, "A comparative analytic study on the gaussian mixture and context dependent deep neural network hidden markov models," in *INTERSPEECH*. ISCA, 2014.
- [4] B. B. Paul and E. A. Martin, "Speaker stress-resistant continuous speech recognition," in *International Conference on Acoustics, Speech and Signal Processing*. IEEE, 1988.
- [5] H. Hermansky, D. P. Ellis, and S. Sharma, "Tandem connectionist feature extraction for conventional hmm systems," in *International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2000.
- [6] M. L. Seltzer, D. Yu, and Y. Wang, "An investigation of deep neural networks for noise robust speech recognition," in *International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2013.
- [7] T. Yoshioka, A. Sehr, M. Delcroix, K. Kinoshita, R. Maas, T. Nakatani, and W. Kellermann, "Making machines understand us in reverberant rooms," in *IEEE Signal Processing Letter*. IEEE, 2012.
- [8] F. Weninger, H. Erdogan, S. Watanabe, E. Vincent, J. L. Roux, J. R. Hershey, and B. Schuller, "Speech enhancement with lstm recurrent neural networks and its application to noise-robust asr," in *International Conference on Latent Variable Analysis and Single Separation*. IEEE, 2015.
- [9] K. C. Sim, A. Narayanan, A. Misra, A. Tripathi, G. Pundak, T. Sainath, P. Haghani, B. Li, and M. Bacchiani, "Domain adaptation using factorized hidden layer for robust automatic speech recognition," in *INTERSPEECH*, 2018.
- [10] H. Tang, W. N. Hsu, F. Grondin, and J. Glass, "A study of enhancement, augmentation, and autoencoder methods for domain adaptation in distant speech recognition," in *INTERSPEECH*, 2018.
- [11] A. Narayanan, A. Misra, K. C. Sim, G. Pundak, A. Tripathi, M. Elfeky, P. Haghani, T. Strohmaier, and M. Bacchiani, "Toward domain-invariant speech recognition via large scale training," in <https://arxiv.org/abs/1808.05312>, 2018.
- [12] Jimmy Ba and Rich Caruana, "Do deep nets really need to be deep?," in *Advances in neural information processing systems*, 2014, pp. 2654–2662.
- [13] Jinyu Li, Rui Zhao, Jui-Ting Huang, and Yifan Gong, "Learning small-size dnn with output-distribution-based criteria," in *Fifteenth annual conference of the international speech communication association*, 2014.
- [14] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *NIPS Deep Learning and Representation Learning Workshop*, 2015.
- [15] Yevgen Chebotar and Austin Waters, "Distilling knowledge from ensembles of neural networks for speech recognition.," in *Interspeech*, 2016, pp. 3439–3443.
- [16] Shan You, Chang Xu, Chao Xu, and Dacheng Tao, "Learning from multiple teacher networks," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2017, pp. 1285–1294.
- [17] Zhongdi Qu, Parisa Haghani, Eugene Weinstein, and Pedro Moreno, "Syllable-based acoustic modeling with ctc-smbr-lstm," in *Automatic Speech Recognition and Understanding Workshop (ASRU), 2017 IEEE*. IEEE, 2017, pp. 173–177.
- [18] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," in *The Journal of the Acoustical Society of America*. IEEE, 1979, p. 943950.
- [19] I. Himawan, P. Motlicek, D. Imseng, B. Potard, N. Kim, and J. Lee, "Learning feature mapping using deep neural network bottleneck features for distant large vocabulary speech recognition," in *International Conference on Acoustics, Speech and Signal Processing*, 2015.
- [20] Shiliang Zhang, Ming Lei, Zhijie Yan, and Lirong Dai, "Deep-fsmn for large vocabulary continuous speech recognition," *arXiv preprint arXiv:1803.05030*, 2018.
- [21] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, 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 EPFL-CONF-192584.
- [22] K. Chen and Q. Huo, "Scalable training of deep learning machines by incremental block training with intra-block parallel optimization and blockwise model-update filtering," in *ICASSP*. IEEE, 2016, p. 58805884.