

Improving DNNs Trained With Non-Native Transcriptions Using Knowledge Distillation and Target Interpolation

Amit Das, Mark Hasegawa-Johnson

University of Illinois, USA

amitdas@illinois.edu, jhasegaw@illinois.edu

Abstract

Often, it is hard to find native transcribers in languages that are under-resourced. However, online Turkers (crowd workers) available in online marketplaces can serve as valuable alternative resources for providing transcriptions in the target language. Since the Turkers neither speak nor have any familiarity with the target language, the transcriptions they provide are noisy and are called Probabilistic Transcriptions (PT). Conventional Deep Neural Network (DNN) training using PTs do not necessarily improve error rates over Gaussian Mixture Models (GMMs). Previously reported results have demonstrated some success by adopting the Multi-Task Learning (MTL) approach. In this study, we report further improvements using Knowledge Distillation (KD) and Target Interpolation (TI) as ways to alleviate transcription errors in PTs. In the KD method, knowledge is transferred from a reasonably well-trained multilingual DNN to the target language DNN trained using PTs. In the TI method, the confidences of the labels provided by PTs are modified using the confidences of the target language DNN. Results show an average improvement in phone error rates (PER) by about 2.12% absolute across Swahili, Amharic, Dinka, and Mandarin. **Index Terms:** knowledge distillation, target interpolation, deep neural networks, under-resourced, cross-lingual speech recognition

1. Introduction

A well-resourced language (WRL) is a language (e.g. English) with an abundance of resources to support development of speech technology. Those resources are usually defined in terms of 100+ hours of speech data, corresponding transcriptions, pronunciation dictionaries, and language models. On the contrary, an under-resourced language (URL) lacks one or more of these resources. The most expensive and time consuming resource is the acquisition of transcriptions due to the difficulty in finding native transcribers.

To circumvent this difficulty, transcriptions can be collected from online non-native crowd workers, or Turkers, who neither speak the target language nor have any familiarity with it. Briefly, a single utterance in some target language L is transcribed by multiple Turkers who do not speak L . Due to this, no single Turker can generate the correct transcription. Instead, a collection of transcriptions from multiple Turkers is constructed for a single utterance in L . This collection can be represented as a confusion network. We will refer to such a network as a *Probabilistic Transcription* (PT) [1]. On the contrary, the correct transcription generated by a native speaker can be represented as a single sequence of labels. We will refer to this sequence as a *Deterministic Transcription* (DT). DTs are simply conventional transcriptions that we frequently encounter while building large corpora ASR systems.

As an example, consider the DT for the word “cat” as

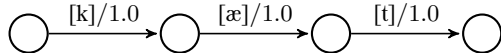


Figure 1: A deterministic transcription (DT) for the word cat.

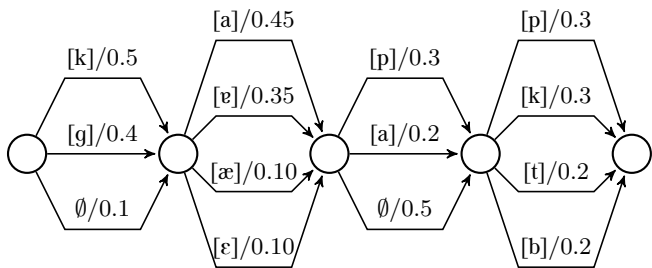


Figure 2: A probabilistic transcription (PT) for the word cat.

shown in Fig. 1. Each arc represents a label and a probability value which is 1.0 always. On the other hand, a PT is the network in Fig. 2. The arc weight specifies the conditional probability that the phoneme was spoken, given the evidence in the transcripts. Because workers cannot distinguish all phone pairs in the utterance language, these weights are usually less than 1.0. In terms of training a DNN, force alignment using DTs result in 1-hot alignments that are frequently observed in conventional transcripts. However, in the case of PTs, the alignments are soft since a single frame could have multiple labels with non-zero probabilities.

Conventional training of DNNs using PTs alone do not necessarily improve error rates over GMMs [2]. Due to this, MTL style training was introduced to train DNNs using a mixture of noisy PTs from the target URL and clean DTs from multiple other WRLs [2–4]. The strong supervision provided by the DTs has the effect of compensating errors caused by PTs.

In this study, we focus on Knowledge Distillation (KD) and Target Interpolation (TI) to further alleviate the effect of noisy labels in PTs. In [5], the authors describe KD as the process of transferring knowledge from a large cumbersome model (or an ensemble of models) to a small distilled model. The cumbersome and distilled model are sometimes referred to as the teacher and student models. Hence, KD is also known as Teacher-Student (TS) learning. If \mathcal{D} is a data set on which the student model is to be trained, then the DNN training procedure involves the following steps. In the first step, feedforward \mathcal{D} through a well-trained teacher DNN to generate the posterior outputs (or teacher labels). The teacher labels form a soft target distribution for each training example in \mathcal{D} . In the second step, train the student DNN by minimizing the cross-entropy (CE) loss between the teacher labels and posterior outputs of

the student DNN. Thus, the student DNN attempts to mimic the behavior of the teacher DNN by trying to match its own outputs with those of the teacher labels. To improve the generalizability of the student DNN, the teacher labels could be generated by using a high temperature T in the softmax of the teacher DNN. The same temperature T is then used at the softmax of the student DNN during CE training. It can be shown that when $T \rightarrow \infty$ (high temperature limit), CE training is equivalent to minimizing the mean square error (MSE) of the logits (pre-softmax activations) between the teacher and student DNNs [5].

Several studies [6–15] in the past have used KD to improve DNNs. In [6], a small DNN was trained using teacher labels generated by feedforwarding a large number of untranscribed data through a large DNN. In other studies, the authors transfer the knowledge from a large RNN to a small DNN [7] or from a large DNN to a small highway DNN [8]. In [9, 10], KD was used to improve robustness of DNNs to noisy data. The one that is most relevant to our work is in [11] where KD was used for adaptation to low resource Japanese dialects and children’s speech.

In the TI approach, we interpolate the confidences of the labels provided by PTs with the confidences of the target language DNN. The DNN is then trained using the new interpolated confidence values. Intuitively, we emphasize the beliefs of the learner rather than completely relying on noisy “ground truth” labels.

The remainder of the paper is organized as follows. In Section 2 and Section 3, we describe the KD and TI frameworks respectively. In Section 4, we discuss our experiments and results. In Section 5, we present our conclusions.

2. Knowledge Distillation (KD)

In this section, we provide a brief outline of the KD framework. Consider an input feature vector \mathbf{x} . A generalized softmax is a softmax function operating on logits $z_k(\mathbf{x})$ and a temperature $T \in \mathbb{R}^+$. Here, $k \in \{1, \dots, K\}$, where K is the total number of labels. We will denote $z_k(\mathbf{x})$ as simply z_k and assume the dependence on \mathbf{x} is implicit. The output of the generalized softmax $y_k(T)$ is given by,

$$y_k(T) = \frac{\exp(z_k/T)}{\sum_{j=1}^K \exp(z_j/T)}. \quad (1)$$

There are two extreme cases in Eq. (1). Let $\mathbf{y}(T) = [y_1(T) \dots y_K(T)]'$. For very hot and cold temperatures, $\mathbf{y}(T)$ approaches the uniform and 1-hot distribution respectively. Thus, $\lim_{T \rightarrow \infty} y_k(T) = \frac{1}{K}$ and $\lim_{T \rightarrow 0} y_k(T) = \mathbb{1}_{[k=\arg \max_{1 \leq j \leq K} y_j]}$. In the KD framework, the student model is trained to minimize the loss,

$$E_{\text{KD}} = \rho C(\mathbf{p}, \mathbf{y}(1)) + (1 - \rho) C(\mathbf{q}(T), \mathbf{y}(T)), \quad (2)$$

where,

$$C(\mathbf{p}, \mathbf{y}(1)) = - \sum_{k=1}^K p_k \log y_k(1), \quad (3)$$

$$C(\mathbf{q}(T), \mathbf{y}(T)) = - \sum_{k=1}^K q_k(T) \log y_k(T). \quad (4)$$

The term p_k is the posterior probability of label k given the feature vector \mathbf{x} . Since this is generated from the PTs, it need not be a binary value 0 or 1. Thus, \mathbf{p} need not be a 1-hot vector.

It is a soft target distribution. Likewise, $q_k(T)$ is the posterior probability of label k generated by feedforwarding \mathbf{x} through a teacher DNN equipped with a generalized softmax with temperature T . In other words, it is a teacher label. In the under-resourced scenario, the teacher DNN is a multilingual DNN trained with DTs from WRLs. The term $y_k(T)$ is the posterior probability of label k generated by feedforwarding \mathbf{x} through a student DNN equipped with a generalized softmax with temperature T as in Eq. (1). The student DNN is a target language DNN to be trained with PTs from the URL. The outputs of the student DNN in Eq. (3) is constrained to a temperature of one whereas in Eq. (4) the temperature can be any $T \in \mathbb{R}^+$. Finally, ρ is a weight that balances the losses in Eq. (3) and Eq. (4).

During backpropagation, the gradient of Eq. (4) with respect to the student logit z_k , i.e., $\frac{\partial C(\mathbf{q}, \mathbf{y})}{\partial z_k}$, is artificially scaled by T^2 . This is because the gradient itself is a function of $1/T^2$. Thus, the artificial scaling removes the dependence on T . As a result, the individual backpropagation errors from Eq. (3) and Eq. (4) have similar scales and can be added meaningfully.

Knowledge distillation specializes to several interesting cases. When $\rho = 1$, Eq. (2) is same as the standard cross-entropy loss. When $0 < \rho < 1$ and $T = 1$, Eq. (2) is equivalent to regularizing the cross-entropy loss with Kullback-Leibler Divergence (KLD) [16]. When $\rho = 0$ (indicating absence of ground truth labels), Eq. (2) can be used for unsupervised adaptation. Thus, for the case of $\rho = 0$, $T = 1$ and when the student DNN is not initialized from a teacher DNN, Eq. (2) can be used for unsupervised adaptation especially if the teacher DNN is a large reliable model [6]. When $\rho = 0$, $T = 1$ and the student DNN is initialized from the teacher DNN, training using Eq. (2) is equivalent to self-training. Here, the teacher labels $\mathbf{q}(1)$ are identical to the outputs $\mathbf{y}(1)$ of the student DNN only before training begins. However, once training begins, the teacher labels are kept constant whereas the student outputs are allowed to change with every weight update.

3. Target Interpolation (TI)

In this section, we provide a brief outline of the TI framework. We will omit the dependence on T since in this section $T = 1$ always. First, we define $C(f(\mathbf{y}), \mathbf{y})$ as,

$$C(f(\mathbf{y}), \mathbf{y}) = - \sum_{k=1}^K f(y_k) \log y_k, \quad (5)$$

where $f(\cdot)$ is an element-wise function of \mathbf{y} . The DNN is trained to minimize the loss,

$$\begin{aligned} E &= \rho C(\mathbf{p}, \mathbf{y}) + (1 - \rho) C(f(\mathbf{y}), \mathbf{y}), \\ &= C(\rho \mathbf{p} + (1 - \rho) f(\mathbf{y}), \mathbf{y}), \end{aligned} \quad (6)$$

where $C(\mathbf{p}, \mathbf{y})$ is as defined in Eq. (3). The second step in Eq.6 is due to the linearity of $C(\cdot, \cdot)$ in the first argument. We consider two among several choices of $f(\cdot)$. They are,

$$f(y_k) = \begin{cases} y_k, & \text{(soft)} \\ \mathbb{1}_{[k=\arg \max_{1 \leq j \leq K} y_j]}, & \text{(hard)} \end{cases} \quad (7)$$

With these choices, the corresponding loss functions are,

$$E_{\text{soft}} = - \sum_{k=1}^K (\rho p_k + (1 - \rho) y_k) \log y_k, \quad (8)$$

$$E_{\text{hard}} = - \sum_{k=1}^K (\rho p_k + (1 - \rho) \mathbb{1}_{[k=\arg \max_{1 \leq j \leq K} y_j]}) \log y_k. \quad (9)$$

And, the error gradients are,

$$\frac{\partial E_{\text{soft}}}{\partial z_k} = \rho(y_k - p_k) + (1 - \rho)y_k(I(y_k) - H(\mathbf{y})), \quad (10)$$

$$\frac{\partial E_{\text{hard}}}{\partial z_k} = \rho(y_k - p_k) + (1 - \rho)(y_k - \mathbb{1}_{[k=\arg \max_{1 \leq j \leq K} y_j]}), \quad (11)$$

where,

$$I(y_k) = - \log y_k,$$

$$H(\mathbf{y}) = - \sum_{k=1}^K y_k \log y_k.$$

The motivation behind the choices in Eq. (7) is that we use the label confidences of the DNN instead of completely relying on the noisy PT labels. Hence, we modify the PT confidence p_k with a new confidence which is an interpolation between p_k and $f(y_k)$. For the soft case, we use the entire output distribution of the DNN. Then the loss in Eq. (8) becomes the standard cross-entropy loss with entropy regularization. A DNN trained using this loss function will find a balance between minimizing the cross-entropy loss $C(\mathbf{p}, \mathbf{y})$ while also lowering the entropy of its outputs $C(\mathbf{y}, \mathbf{y})$. Since PTs are prone to high entropies, lowering the entropies of the DNN outputs is desirable. For the hard case, we simply binarize the DNN outputs to a 1-hot distribution. Compared to the soft case, the hard case ignores the cross-correlation between different classes. In both cases, the new interpolated confidences still form a valid probability distribution since they sum to one when summed over the labels.

4. Experiments and Results

4.1. Data

Multilingual audio files were obtained from the Special Broadcasting Service (SBS) network which publishes multilingual radio podcasts in Australia. The corpus is summarized in Table 1. Natively transcribed DTs in Arabic (*arb*), Cantonese (*yue*), and Hungarian (*hun*) were always treated as data from source WRLs. PTs from Turkers were used as training data for the target URL. We experimented with four target URLs - Swahili (*swh*), Amharic (*amh*), Dinka (*din*), and Mandarin (*cmn*) - in a round-robin fashion. For example, if *swh* is the target language, then the training set consists of PTs in *swh* and DTs in the remaining six languages (*amh*, *din*, *cmn*, *arb*, *yue*, *hun*). Thus, it excludes having *swh* DTs in the training set. In this sense, our experiments fall under the domain of zero-resource speech recognition.

More than 2500 Turkers participated in transcribing, with roughly 30% of them claiming to know only English. The remaining Turkers claimed knowing other languages such as Spanish, French, German, Japanese, and Mandarin. The utterances were limited to a length of 5 seconds. This is because the Turkers did not understand the utterance language and it was easier for them to annotate short utterances than long. Since English was the most common language among the Turkers,

Table 1: *SBS Multilingual Corpus*.

Language	Utterances		Phones
	Train	Test	
Swahili (<i>swh</i>)	462	123	48
Amharic (<i>amh</i>)	516	127	37
Dinka (<i>din</i>)	248	53	27
Mandarin (<i>cmn</i>)	467	113	52
Arabic (<i>arb</i>)	468	112	46
Cantonese (<i>yue</i>)	544	148	32
Hungarian (<i>hun</i>)	459	117	65
All	-	-	82

they were asked to annotate the sounds using English letters. The sequence of letters was not meant to be meaningful English words or sentences since this would be detrimental to the final performance. The important criterion was that the annotated letters represent sounds they heard from the utterances as if they were listening to a sequence of nonsense syllables in some exotic language. Since no Turker is likely to generate the perfect transcript, each utterance was transcribed by ten Turkers creating ten different transcripts per utterance. These transcripts were converted to phones and merged into a PT using [1]. Turkers were typically paid \$500 per ten Turkers for transcribing an hour of audio. As for DTs, the same set of utterances were transcribed by native speakers in the target language. However, the DTs in the target language were used only to know the ground truth hypotheses which are necessary for evaluating the ASR performance on the test set.

The training set consists of a) about 40 minutes of PTs in the target URL and, b) about 40 minutes of DTs in multiple WRLs. The development and test sets were worth 10 minutes each. The test utterances were randomly selected to avoid any speaker or gender bias. Going back to our previous example, if *swh* is the target language, then the training set consists of 40 minutes of PTs in *swh* and 40 minutes of DTs each in *amh*, *din*, *cmn*, *arb*, *yue*, *hun* (total $40 \times 6 = 240$).

All experiments were conducted using the Kaldi toolkit [17]. Kaldi source code in C++ and toy examples under the proposed KD and TI frameworks are available in our github repository.¹

4.2. Experiments

In this section, we describe the features, baseline, and the proposed experiments. Thirteen Mel Frequency Cepstral Coefficients (MFCCs), spliced with +/- 3 neighboring frames, were extracted from speech utterances. These were then transformed using a Linear Discriminant Analysis (LDA) transform followed by Feature-Space Maximum Likelihood Linear Regression (fMLLR) transform resulting in 40-dimensional fMLLR features. These features were kept low dimensional to avoid the curse-of-dimensionality problem which is more likely for under-resourced scenarios. These features were then mean normalized using Cepstral Mean Normalization (CMN) before using them for DNN training.

As for the labels in DNN training, the forced aligned senones obtained from HMM models were treated as the ground truth labels for DNN training. Since a PT is a lattice, forced alignment performed on a PT produces an *alignment lattice* instead of an alignment sequence. Running forward-backward recursion on the alignment lattice generates the frame-level poste-

¹`git clone -b teacher-student`
<https://github.com/irrawaddy28/SBS-kaldi-2015>

Table 2: Phone error rates of different MTL systems trained with CE, KLD, and KD losses. The parameters ρ and T are the weighting and temperature parameters in Eq. (2).

System	Parameters		Language			
	ρ	T	swh	amh	din	cmn
Baseline (CE)	1	-	44.89	60.79	58.65	53.53
KLD	0.6	1	44.11	59.97	58.19	51.00
KLD	0.4	1	44.21	59.36	58.33	50.29
KLD	0.2	1	44.63	59.55	58.65	50.93
KD	0.6	2	44.12	59.82	58.15	50.93
KD	0.4	2	43.66	59.40	57.97	49.85
KD	0.2	2	44.40	59.08	58.26	49.38

rriors which are soft as opposed to 1-hot. Phone based language models (LMs) were built from text data in the target language mined from Wikipedia. Consequently, a phone based decoder was used to generate the final ASR hypotheses. These were evaluated using PERs. Results from the following experiments were compared for evaluation:

- **Baseline [4]:** An MTL system was trained consisting of six shared hidden layers and two separate softmax layers (one softmax per task). The shared hidden layers of the MTL system were initialized from a multilingual DNN. Both the tasks were trained to minimize the CE loss. However, the targets at the first softmax are PTs in a target URL. The targets at the second softmax are DTs in the remaining six WRLs. We do not train with DTs in the target URL.
- **KLD Regularization [16]:** Instead of minimizing the standard CE loss for the first task, here the MTL system was trained to minimize the loss in Eq. (2) for the special case of $T = 1$ and $0 < \rho < 1$. Specifically, $\rho \in \{0.2, 0.4, 0.6, 0.8\}$ was used.
- **Knowledge Distillation:** The first task of the MTL system was trained to minimize the loss E_{KD} in Eq. (2) with $0 < \rho < 1$ and $T > 1$. Specifically, $T \in \{2, 3\}$ and $\rho \in \{0.2, 0.4, 0.6, 0.8\}$ were used.
- **Target Interpolation:** The first task of the MTL system was trained to minimize the loss E_{soft} in Eq. (8) and E_{hard} in Eq. (9) while using $\rho \in \{0.2, 0.4, 0.6, 0.8\}$.

4.3. Results

The PERs comparing the MTL systems trained with CE, KLD, and KD losses are outlined in Table 2. The systems are named eponymously after their loss types. We highlight only the most interesting cases with ρ in the range 0.6 – 0.2 and $T = 2$. From Table 2, it is clear that the KD systems outperform the baseline CE and KLD systems.

Now, we analyze the effect of T and ρ on recognition rates. Keeping ρ fixed and varying T is equivalent to comparing KLD with KD systems. Thus, as T increases (keeping ρ constant), KD systems outperform their KLD counterparts most of the times. Increasing T makes the class correlations more pronounced. This indicates that the temperature parameter improves the generalization capacity of the DNNs by avoiding tuning to the noisy PTs. Next, keeping $T > 1$ fixed and varying ρ is equivalent to comparing within the family of KD systems. As ρ decreases, the PERs tend to decrease first and then increase. Desirable values of ρ are $\rho < 0.5$. This implies that the performance improves when the system relies increasingly on the teacher labels rather than the PT labels. However, this trend reverses for very low values of ρ . For example, in the extreme case where $\rho = 0$ (ignore PT labels), we noticed exceedingly

Table 3: Phone error rates of different MTL systems trained with CE and TI losses. The parameter ρ is the weighting parameter in Eq. (8) and Eq. (9).

System	Parameter	Language			
	ρ	swh	amh	din	cmn
Baseline (CE)	1.0	44.89	60.79	58.65	53.53
TI (Hard)	0.6	43.96	60.44	58.69	51.14
TI (Hard)	0.4	44.08	59.98	57.94	49.81
TI (Hard)	0.2	44.24	60.58	59.19	51.20
TI (Soft)	0.6	43.49	60.19	58.62	51.09
TI (Soft)	0.4	43.29	59.65	57.65	50.02
TI (Soft)	0.2	44.16	61.14	59.26	50.79

Table 4: Summary of the best proposed systems. Absolute improvements over the baseline system inside parantheses.

Lang	Baseline (CE)	Best		Parameters	
	PER	PER	System	ρ	T
swh	44.89	43.29 (1.60)	TI (Soft)	0.4	-
amh	60.79	59.08 (1.71)	KD	0.2	2
din	58.65	57.65 (1.00)	TI (Soft)	0.4	-
cmn	53.53	49.38 (4.15)	KD	0.2	2

high PERs above 85%. This proves that PT labels are still useful.

The PERs comparing the CE and TI systems are outlined in Table 3. Again, we highlight only the most interesting cases of ρ (0.6 – 0.2). Clearly, both variants of TI systems outperform the baseline CE system. Among the TI systems, TI (Soft) outperforms TI (Hard) for the African languages (Swahili, Amharic, and Dinka). For Mandarin, TI (Hard) outperforms TI (Soft) by a small margin. Surprisingly, for both TI (Hard) and TI (Soft), $\rho = 0.4$ is the most desirable value. Moreover, quite conveniently, this value of ρ does not change across languages explored in this study. Similar to the KD system, values of $\rho < 0.5$ imply that the performance improves when the system relies increasingly on the DNN labels rather than the PT labels. This means that the new interpolated targets are effective in alleviating the noise in PT labels. However, similar to the KD system, setting $\rho = 0$ results in very high PERs.

Finally, a summary of the best proposed systems for each language, along with their parameters, is highlighted in Table 4. The average improvement is about 2.12% absolute. This is quite useful for us considering that this is a zero-resource scenario and we do not have access to reliable ground truth DTs in the target URL.

We conducted additional experiments in an attempt to further boost the performance of the best KD systems. Since the PT distribution \mathbf{p} is a soft distribution, we parameterized \mathbf{p} with a new temperature parameter T_{PT} . After changing \mathbf{p} to $\mathbf{p}(T_{PT})$, we minimize the KD loss E_{KD} in Eq.(2). We noticed an improvement in PER over the best KD systems by about 0.2% absolute when $T_{PT} = 2$. Since the improvement is marginal, we continue to investigate ways to improve \mathbf{p} .

5. Conclusions

In this study, we report further improvements in DNNs trained with noisy non-native transcriptions (PTs) while not having access to native transcriptions (DTs) in the target language. We proposed Knowledge Distillation and Target Interpolation to alleviate the effect of noise in PTs. We reported consistent improvements in recognition rates for all languages explored in this study with an average improvement of 2.12% absolute.

6. References

- [1] P. Jyothi and M. Hasegawa-Johnson, "Transcribing Continuous Speech Using Mismatched Crowdsourcing," in *Interspeech*, 2015, pp. 2774–2778.
- [2] A. Das and M. Hasegawa-Johnson, "An Investigation on Training Deep Neural Networks Using Probabilistic Transcriptions," in *Interspeech*, 2016, pp. 3858–3862.
- [3] V. H. Do, N. F. Chean, B. P. Lim, and M. Hasegawa-Johnson, "Multi-Task Learning Using Mismatched Transcription for Under-Resourced Speech Recognition," in *Interspeech*, 2017, pp. 2073–2077.
- [4] A. Das, M. Hasegawa-Johnson, and K. Veselý, "Deep Autoencoder Based Multi-Task Learning Using Probabilistic Transcriptions," in *Interspeech*, 2017, pp. 2073–2077.
- [5] G. Hinton, O. Vinyals, and J. Dean, "Distilling the Knowledge in a Neural Network," in *arXiv:1503.02531*, 2015.
- [6] J. Li, R. Zhao, J.-T. Huang, and Y. Gong, "Learning Small-Size DNN with Output-Distribution-Based Criteria," in *Interspeech*, 2014.
- [7] W. Chan, N. R. Ke, and I. Lane, "Transferring Knowledge from a RNN to DNN," in *Proc. Interspeech*, 2015, pp. 3264–3268.
- [8] L. Lu, M. Guo, and S. Renals, "Knowledge Distillation for Small-Footprint Highway Networks," in *Proc. ICASSP*, 2017, pp. 4820–4824.
- [9] K. Markov and T. Matsui, "Robust Speech Recognition Using Generalized Distillation Framework," in *Interspeech*, 2016, pp. 2364–2368.
- [10] S. Watanabe, T. Hori, J. L. Roux, and J. Hershey, "Student-Teacher Network Learning with Enhanced Features," in *ICASSP*, 2017.
- [11] T. Asami, R. Masumura, Y. Yamaguchi, H. Masataki, and Y. Aono, "Domain Adaptation of DNN Acoustic Models Using Knowledge Distillation," in *Proc. ICASSP*, 2017.
- [12] J. Cui, B. Kingsbury, B. Ramabhadran, G. Saon, T. Sercu, K. Audhkhasi, A. Sethy, M. Nussbaum-Thom, and A. Rosenberg, "Knowledge Distillation Across Ensembles of Multilingual Models for Low-Resource Languages," in *Proc. ICASSP*, 2017, pp. 4825–4829.
- [13] Y. Chebotar and A. Waters, "Distilling Knowledge from Ensembles of Neural Networks for Speech Recognition," in *Proc. Interspeech*, 2016, p. 34393443.
- [14] G. Tucker, M. Wu, M. Sun, S. Panchapagesan, G. Fu, and S. Vitaladevuni, "Model Compression Applied to Small-Footprint Keyword Spotting," in *Proc. Interspeech*, 2016, p. 18781882.
- [15] J. Li, R. Zhao *et al.*, "Developing Far-Field Speaker System via Teacher-Student Learning," in *Proc. ICASSP*, 2018.
- [16] D. Yu, K. Yao, H. Su, G. Li, and F. Seide, "KL-Divergence Regularized Deep Neural Network Adaptation For Improved Large Vocabulary Speech Recognition," in *Proc. ICASSP*, 2013, pp. 7893–7897.
- [17] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlíček, Y. Qian, P. Schwarz, J. Silovský, G. Stemmer, and K. Veselý, "The Kaldi Speech Recognition Toolkit," 2011.