# JOINT MODELING OF ACCENTS AND ACOUSTICS FOR MULTI-ACCENT SPEECH RECOGNITION

*Xuesong Yang*[*], *Kartik Audhkhasi*[†], *Andrew Rosenberg*[†], *Samuel Thomas*[†],
*Bhuvana Ramabhadran*[†], *Mark Hasegawa-Johnson*[*]

[*]University of Illinois at Urbana-Champaign, Urbana, IL
[†]IBM T. J. Watson Research Center, Yorktown Heights, NY

## ABSTRACT

The performance of automatic speech recognition systems degrades with increasing mismatch between the training and testing scenarios. Differences in speaker accents are a significant source of such mismatch. The traditional approach to deal with multiple accents involves pooling data from several accents during training and building a single model in multi-task fashion, where tasks correspond to individual accents. In this paper, we explore an alternate model where we jointly learn an accent classifier and a multi-task acoustic model. Experiments on the American English Wall Street Journal and British English Cambridge corpora demonstrate that our joint model outperforms the strong multi-task acoustic model baseline. We obtain a 5.94% relative improvement in word error rate on British English, and 9.47% relative improvement on American English. This illustrates that jointly modeling with accent information improves acoustic model performance.

***Index Terms***— End-to-end models, acoustic modeling, multi-accent speech recognition, multi-task learning

## 1. INTRODUCTION

Recent breakthroughs in automatic speech recognition (ASR) have resulted in a word error rate (WER) on par with human transcribers [1, 2] on the English Switchboard benchmark. However, dealing with acoustic condition mismatch between the training and testing data is a significant challenge that still remains unsolved. It is well-known that the performance of ASR systems degrades significantly when presented with speech from speakers with different accents, dialects and speaking styles than those encountered during system training [3]. In this paper, we specifically focus on acoustic modeling for multi-accent ASR.

Dialects are defined as variations within a language that differ in geographical regions and social groups, which can be distinguished by traits of phonology, grammar, and vocabulary [4]. Specifically, dialects may be associated with the residence, ethnicity, social class, and native language of speakers. For example, in British and American English, same words can have different spellings, like *favour* and *favor*; or different pronunciations, such as ˈʃɛdjuːl in UK English vs. ˈskɛdʒʊl in US English for the word *schedule*; in Spanish, vocabulary may evolve differently between dialects, like for the phrase *cell phone*, Castilian Spanish uses *móvil* while Latin American use *celular* [5]; in English, same phoneme may be realized differently, phoneme /e/ in *dress* is pronounced as /ɛ/ in England and /e/ in Wales; in Arabic, dialects may also differ in intonation and rhythm cues [6]. In this paper, we focus on the issue of differing pronunciations, while eschewing considerations of grammatical and vocabulary differences.

Acoustic modeling across multiple accents has been explored for many years, and various approaches can be summarized into three categories - *Unified models*, *Adaptive models*, and *Ensemble models*. A unified model is trained on a limited number of accents, and can be generalized to any accent [7, 8]. An adaptive model fine-tunes the unified model on accent-specific data assuming that the accent is known [9–11]. An ensemble model aggregates all accent-specific recognizers, and produces an optimal model by selection or combination for recognition [5, 12, 13]. Experiments have revealed that the unified model usually underperforms the adaptive model, which in turn underperforms the ensemble model [7, 8].

We note that these prior approaches do not explicitly include accent information during training, but do so only indirectly, for example, through the different target phoneme sets for various accents. This contrasts sharply with the way in which humans memorize the phonological and phonetic forms of accented speech: "mental representations of phonological forms are extremely detailed," and include "traces of individual voices or types of voices" [14]. In this paper, we propose to link the training of ASR acoustic models and accent identification models, in a manner similar to the linking of these two learning processes in human speech perception. We show that this joint model not only performs well on ASR, but also on accent identification when compared to separately-trained models. Given the recent success in end-to-end models [15–25], we use a bidirectional long short-term memory (BLSTM) recurrent neural network (RNN) acoustic model trained with the connectionist temporal classification (CTC) loss function for acoustic modeling. The accent identification (AID) network is also a BLSTM, but includes an average pooling layer to compute an utterance-level accent embedding. We also introduce a joint architecture where the lower layers of the network are trained using AID as the auxiliary task while multi-accent acoustic modeling remains the primary task of the network.

Next, we use the AID network as a hard switch between the accent-specific output layers of the CTC AM. Preliminary experiments on the Wall Street Journal American English and Cambridge British English corpora demonstrate that our joint model with the AID-based hard-switch achieves lower WER when compared with the state-of-the-art multi-task AM. We also show that the AID model also benefits from joint training.

The remainder of this paper is organized as follows: Section 2 reviews relevant literature. Section 3 introduces our AID model, multi-accent acoustic model, and switching strategy. Section 4 shows experiments and analysis, followed by the conclusion in Section 5.

## 2. RELATED WORK

The most closely related work to ours is from [8], which illustrated that hierarchical grapheme-based AM with auxiliary phoneme-based AMs in four English dialects trained with CTC significantly outperformed accent-specific AMs and grapheme-based AM, respectively, while achieving competitive WER with phoneme-based multi-accent AM. Similarly, Yi et al [11] also trained a multi-accent phoneme-based AM with CTC loss, but instead, adapted accent-specific output layer using its target accent.

Other relevant work compared the performance of training accent or dialect specific acoustic models and joint models. These approaches predicted context-dependent (CD) triphone states using DNNs, and used a weighted finite state transducer (WFST)-based decoder. For example, senones on accents of Chinese are predicted by assuming all accents within a language share a common CD state inventory [9, 10]. Elfeky et al [7] implemented a dialectal multi-task learning (DMTL) framework on three dialects of Arabic using the prediction of a unified set of CD states across all dialects prediction as the primary task and dialect identification as the secondary task. DMTL model deviated from ours in that it directly predicted CD states using convolutional-LSTM-DNNs (CLDNN), and was trained with either cross-entropy or state-level minimum Bayes risk, while ignoring the secondary dialect identification output at recognition time. This DMTL model was trained on all dialectal data and underperformed the dialect-specific model. Dialectal knowledge distilled (DKD) model was also designed in [7], which achieved results competitive to, but below, dialect-specific models.

The effectiveness of dialect-specific models motivated investigations into how to use ensemble methods on multiple dialect-specific acoustic models for recognition. Soto et al [5] explored approaches of selecting and combining the best decoded hypothesis from a pool of dialectal recognizers. This work is still different from ours in that we make our selection directly using predicted dialect. Huang et al [3] used a similar strategy to ours by identifying accent first followed by acoustic model selection, however, this work only considered GMMs as the classifier.

## 3. METHOD

Our proposed system consists of multiple accent-specific acoustic models and accent identification model. We will describe these components and their joint model in this section. Acoustic model selection based on the hard-switch between accent-specific models is illustrated in Section 3.4.
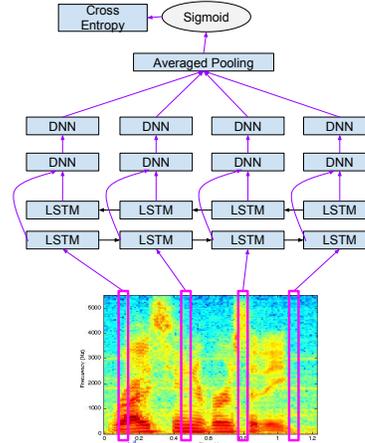
### 3.1. Accent Identification

Accurate identification of a speaker's accent is essential to the pipelined ASR systems, since accent identification (AID) errors can cause large mismatch to acoustic models. Given the hypothesis that accents can be discriminated by spectral features, researchers have attempted to model the spectral distribution of each accent using GMMs. Recently, DNNs have been explored as a much more expressive model compared to GMMs, especially in modeling probability distributions.

We implemented an independent AID that summarizes low-level acoustic features of an utterance by a stack of bidirectional LSTMs (BLSTMs) and DNN projection layers. An average-pooling layer is applied on top of transformed acoustic features, because the acoustic realization of a speaker's accent may not be observable in each frame. Applying average-pooling gives us a more robust estimate of

accent-dependent acoustic features. We note that we assume that the speaker's accent is fixed over the entire utterance.

Figure 1 depicts details of this AID model. A single sigmoidal neuron is used at the output layer for classification because we are only classifying between accents of English - US and UK. We trained the AID network using the cross-entropy loss function.



**Fig. 1**: Proposed accent identification (AID) model with BLSTMs and average-pooling.

### 3.2. Multi-Accent Acoustic Modeling

Recently, end-to-end (E2E) systems have achieved comparable performance to traditional pipelined systems such as hybrid DNN-HMM systems. These E2E systems come with the benefit of avoiding time-consuming iterations between alignment and model building. RNNs using the CTC loss function are a popular approach to E2E systems [15]. The CTC loss computes the total likelihood of the output label sequence given the input acoustics over all possible alignments. It achieves this by introducing a special *blank* symbol that augments the label sequence to make its length equal to the length of the input sequence. Clearly, there are multiple such augmented sequences, and CTC uses the forward-backward algorithms to efficiently sum the likelihoods of such sequences. The CTC loss is

$$p(\mathbf{l}|\mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\mathbf{l})} p(\pi|\mathbf{x}) \tag{1}$$

where $\mathbf{l}$ is the output label sequence, $\mathbf{x}$ is the input acoustic sequence, $\pi$ is a blank-augmented sequence for $\mathbf{l}$, and $\mathcal{B}^{-1}(\mathbf{l})$ is the set of all such sequences. During decoding, the target label sequences can be obtained by either greedy search or a WFST-based decoder.

Our multi-accent acoustic model combines two CTC-based AMs, one for each accent. We applied multiple BLSTM layers shared by two accents to capture accent-independent acoustic features, and placed separate DNNs for each AM to extract accent-specific features. Figure 2 describes the structure of multi-accent acoustic model. Both AMs are jointly trained with an average of the two accent-specific CTC losses.

At test time, this multi-accent model requires knowledge of the speaker's accent to pick out of the two accent-specific targets. We experimented with both the oracle accent label, and using a trained AID network to make this decision.
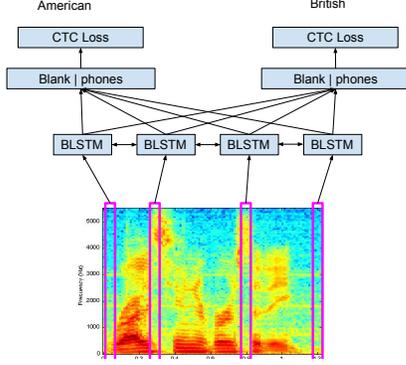
**Fig. 2**: This figure shows the multi-accent acoustic model.

### 3.3. Joint Acoustic Modeling with AID

The previous multi-accent model assumes that multi-tasking between the phone sets of the two accents is sufficient to make the network learn accent-specific acoustic information. An alternate approach is to explicitly supervise the network with accent information. This leads us to our joint model, with multi-accent acoustic modeling as primary tasks at higher layers, and with AID as an auxiliary task at lower level layers, as shown in Figure 3. This joint model aggregates two modules with the same structures to the forementioned models in Section 3.1 and 3.2, and can be jointly trained in an end-to-end fashion with the objective function,

$$\min_{\Theta} \mathcal{L}_{\text{Joint}}(\Theta) = (1 - \alpha) * \mathcal{L}_{\text{AM}}(\Theta) + \alpha * \mathcal{L}_{\text{AID}}(\Theta)$$

where $\alpha$ is an interpolation weight balancing between CTC loss of multi-accent AMs and the cross-entropy loss of AID, and $\Theta$ is the model parameters. CTC loss $\mathcal{L}_{\text{AM}}$ sums up the probabilities of all possible paths corresponding to Equation (1), while AID classification loss $\mathcal{L}_{\text{AID}}$ is cross-entropy. The two losses are at different scales, so the optimal value of $\alpha$ needs to be tuned on development data.
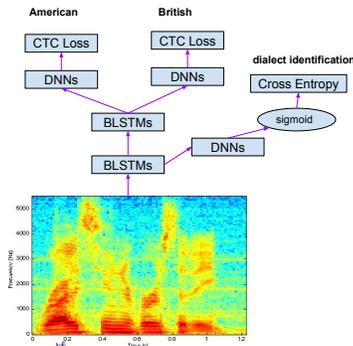


**Fig. 3**: Proposed joint model for accent identification and acoustic modeling.

### 3.4. Model Selection by Hard-Switch

Given a trained CTC-based multi-accent acoustic model and AID classifier, we apply maximum likelihood estimation to switch between the accent-specific output layers $\mathbf{y}_{\text{US}}$ and $\mathbf{y}_{\text{UK}}$. Let $P_{\text{AID}}(\text{US}|\mathbf{x})$ denote the probability of the US accent estimated by AID. We

threshold this probability at $0.5$ to obtain the accent hard-switch $s_{\text{AID}}(\text{US}|\mathbf{x})$. Hence, we pick the output layer as follows:

$$\mathbf{y} = \begin{cases} \mathbf{y}_{\text{US}} & \text{if } s_{\text{AID}}(\text{US}|\mathbf{x}) = 1 \\ \mathbf{y}_{\text{UK}} & else \end{cases}$$

We note that this strategy applies to both the multi-accent model and the joint model.

### 4. EXPERIMENTS

We perform experiments on two dialects of English corpora–Wall Street Journal-1 American English and Cambridge British English. They contain overlapping, but distinct phone sets of 42 and 45 phones respectively. Both corpora contain approximately 15 hours each of audio. We held-out 5% of the training data as a development set. The window size of each speech frame is 25ms with a frame shift of 10ms. We extracted 40-dimensional log-Mel scale filter banks and performed per-utterance cepstral mean subtraction. We did not use any vocal tract length normalization. We then stacked neighboring frames and picked every alternate frame to get a 80-dimensional acoustic feature stream at half the frame rate. Various models are compared in terms of phone error rate (PER) and word error rate (WER). Particularly, we obtain the PER after simple frame-wise greedy decoding from the DNN projection outputs after removing repeated phones and the *blank* symbol. The Attila toolkit [26] is used to report WER by applying WFST-based decoding. Evaluation is performed on `eval93`[1] American English and `si_dt5b`[2] British English.

Our joint model uses four BLSTM layers where the lowest layer is attached to the AID network and the highest single layer connects to two accent-specific softmax layers. A single DNN layer with 320 hidden units is used for each task. The weights for all models are initialized uniformly from $[-0.01, 0.01]$. Adam [27] optimizer with initial learning rate $5e - 4$ is used, and the gradients are clipped to the range $[-10, 10]$. We discard the training utterances that are longer than 2000 frames. New-bob annealing [28] on the held-out data is used for early stopping, where the learning rate is cut in half whenever the held-out loss does not decrease. For the purpose of fair comparison, we used a four layer BLSTM for the baseline acoustic models as well.

Various models are briefly described as follows:

- ASpec: phoneme-based accent-specific AMs that are trained separately on mono-accent data.

- MTLP: phoneme-based multi-accent AMs that are jointly trained on two accents.

- Joint: proposed phoneme-based joint acoustic model with AID.

### 4.1. Empirical weights for balancing different losses

Our joint model is sensitive to the interpolation weight $\alpha$ between the AM CTC and AID cross-entropy losses. We tuned $\alpha$ on development data. Figure 4 depicts relationship between overall PER of two accents and different $\alpha$ values. When $\alpha$ goes larger, overall PER increases but with small fluctuations, especially at $\alpha$ of 0.01 and 0.2. The PER tends to be the largest if $\alpha$ is 1.0, which is expected since the weights of neural networks are updated only using the AID errors. We found the optimal value of $\alpha$ to be 0.001, which
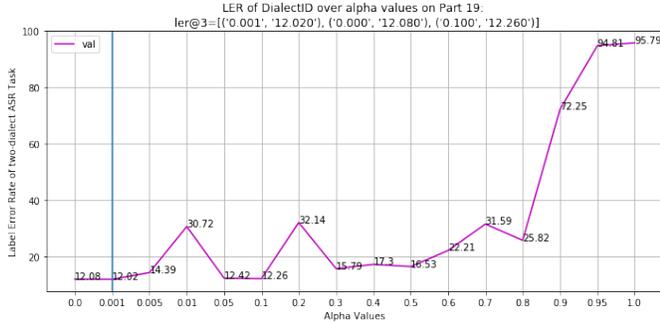
---

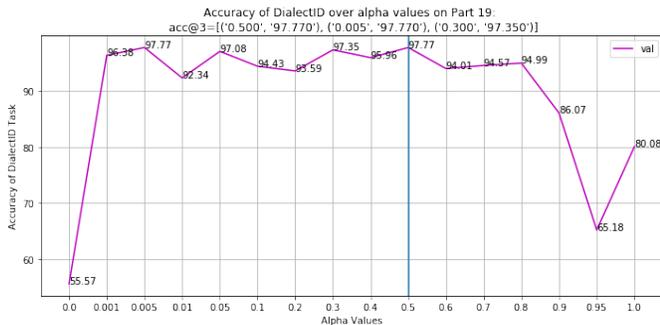**Fig. 4**: PER of joint acoustic model over AID loss weights $\alpha$



**Fig. 5**: Accuracy of AID over various AID weights $\alpha$

achieved minimum PER of 12.02%. Figure 5 illustrates the trend of AID accuracy over different $\alpha$ values. Weights between 0.001 and 0.8 all perform well with accuracies greater than 92%, while tail values lead to even worse performance. When $\alpha$ is 0.5 and 0.005, the best performance is achieved with 97.77% accuracy.

### 4.2. Oracle performance for multi-accent acoustic models

We first evaluate the oracle performance of various models in Table 1. These results assume that the correct accent of each utterance is provided for all models. In other words, the acoustic model corresponds to the correct accent, i.e. the relevant target accent-specific softmax layer is used. It can be seen that the proposed joint model significantly outperforms the accent-specific model (ASpec) by 17.97% relative improvement in overall WER, and multi-task accent model (MTLP) by 6.81%. This observation indicates that deep BLSTM layers shared with multiple accent AMs can learn expressive accent-independent features that refine accent-specific AMs. The auxiliary task, accent identification, also helps by introducing extra accent-specific information. The advantage of augmenting general acoustic features with specific information both implicitly learned by our joint model is observed in natural language processing [29] tasks as well. The value of implicit feature augmentation is a rich area for future investigation.

### 4.3. Hard-switch using distorted AID

The oracle experiments in Section 4.2 demonstrate the value of our proposed joint model and the MTLP model when the AID classifier operates perfectly. This section demonstrates the impact of imperfect AID on the performance using hard-switch. Table 2 shows the

**Table 1**: Oracle performance in word error rates that assumes that the true accent ID is known in advance. Word error rates is calculated after decoding with a WFST-graph incorporating a LM; the relative improvement (*rel.*) for each model over ASpec are reported in the parenthesis.

| corpus | ASpec | MTLP (rel.) | Proposed Model (rel.) |
|---|---|---|---|
| British | 11.5 | 10.1 (-12.17) | 9.5 (-17.39) |
| American | 10.2 | 9.0 (-11.76) | 8.3 (-18.63) |
| average | 10.85 | 9.55 (-11.98) | 8.9 (-17.97) |

results. Given a well-trained independent AID (ind. AID), our joint model still significantly outperforms the two baseline models, and MTLP achieves better WER than ASpec. In comparison to oracle WERs of all models, British WERs are relatively constant without any distortion, however, American English WERs deteriorate accordingly. This is because independent AID has 100% recall for British English utterances on the test data.

It is interesting to note that the biggest improvement over ASpec in WER comes when using the joint model (21.62%) instead of the MTLP model (14.41%) with an independent AID model. The improvement upon further using the AID from the joint model itself is still larger (22.52%). This indicates that the joint model has already learned sufficient accent-specific information through the accent supervision in the lower layers.

**Table 2**: WERs of hard-switch using distorted AID. The *rel.* shows the relative improvement over ASpec; *ind. AID* applies an independent neural AID trained separately. Our *Proposed Model* applies the AID jointly learn with multi-accent AMs.

| Corpus | Pipelines with ind. AID | | | Proposed |
|---|---|---|---|---|
| | ASpec | MTLP (rel.) | Joint (rel.) | Model (rel.) |
| British | 11.5 | 10.1 (-12.17) | 9.5 (-17.39) | 9.5 (-17.39) |
| American | 11.1 | 9.5 (-14.41) | 8.7 (-21.62) | 8.6 (-22.52) |

## 5. CONCLUSION

This paper studies state-of-the-art approaches of acoustic modeling across multiple accents. We note that these prior approaches do not explicitly include accent information during training, but do so only indirectly, for example through the different phone inventories for various accents. We propose an end-to-end multi-accent acoustic modeling approach that can be jointly trained with accent identification. We use BLSTM-RNNs to design acoustic models that can be trained with CTC, and apply an average pooling to compute utterance-level accent embedding. Experiments show that both multi-accent acoustic models and accent identification benefit each other, and our joint model using hard-switch outperforms the state-of-the-art multi-accent acoustic model baseline with a separately-trained AID network. We obtain a 5.94% relative improvement in WER on British English, and 9.47% on American English.

## 6. REFERENCES

[1] Wayne Xiong, Jasha Droppo, Xuedong Huang, Frank Seide, Mike Seltzer, Andreas Stolcke, Dong Yu, and Geoffrey Zweig, "Achieving human parity in conversational speech recognition," *arXiv preprint arXiv:1610.05256*, 2016.

[2] George Saon, Gakuto Kurata, Tom Sercu, Kartik Audhkhasi, Samuel Thomas, Dimitrios Dimitriadis, Xiaodong Cui, Bhuvana Ramabhadran, Michael Picheny, Lynn-Li Lim, et al., "English conversational telephone speech recognition by humans and machines," *arXiv preprint arXiv:1703.02136*, 2017.

[3] Chao Huang, Tao Chen, and Eric Chang, "Accent issues in large vocabulary continuous speech recognition," *International Journal of Speech Technology*, vol. 7, no. 2, pp. 141–153, 2004.

[4] Janet Holmes, *An introduction to sociolinguistics*, Routledge, 2013.

[5] Victor Soto, Olivier Siohan, Mohamed Elfeky, and Pedro Moreno, "Selection and combination of hypotheses for dialectal speech recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 5845–5849.

[6] Fadi Biadsy and Julia Hirschberg, "Using prosody and phonotactics in arabic dialect identification," in *Tenth Annual Conference of the International Speech Communication Association*, 2009.

[7] Mohamed Elfeky, Meysam Bastani, Xavier Velez, Pedro Moreno, and Austin Waters, "Towards acoustic model unification across dialects," in *Spoken Language Technology Workshop (SLT), 2016 IEEE*. IEEE, 2016, pp. 624–628.

[8] Kanishka Rao and Haşim Sak, "Multi-accent speech recognition with hierarchical grapheme based models," in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 4815–4819.

[9] Yan Huang, Dong Yu, Chaojun Liu, and Yifan Gong, "Multi-accent deep neural network acoustic model with accent-specific top layer using the kld-regularized model adaptation," in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.

[10] Mingming Chen, Zhanlei Yang, Jizhong Liang, Yanpeng Li, and Wenju Liu, "Improving deep neural networks based multi-accent mandarin speech recognition using i-vectors and accent-specific top layer," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

[11] Jiangyan Yi, Hao Ni, Zhengqi Wen, Bin Liu, and Jianhua Tao, "Ctc regularized model adaptation for improving lstm rnn based multi-accent mandarin speech recognition," in *Chinese Spoken Language Processing (ISCSLP), 2016 10th International Symposium on*. IEEE, 2016, pp. 1–5.

[12] Yanli Zheng, Richard Sproat, Liang Gu, Izhak Shafran, Haolang Zhou, Yi Su, Daniel Jurafsky, Rebecca Starr, and Su-Youn Yoon, "Accent detection and speech recognition for shanghai-accented mandarin.," in *Interspeech*, 2005, pp. 217–220.

[13] Mohamed Elfeky, Pedro Moreno, and Victor Soto, "Multi-dialectical languages effect on speech recognition: Too much choice can hurt," in *International Conference on Natural Language and Speech Processing (ICNLSP)*, 2015.

[14] Janet Pierrehumbert, "Phonological representation: Beyond abstract versus episodic," *Annu. Rev. Linguist.*, vol. 2, pp. 33–52, 2016.

[15] Alex Graves and Navdeep Jaitly, "Towards end-to-end speech recognition with recurrent neural networks," in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, 2014, pp. 1764–1772.

[16] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al., "Deep speech 2: End-to-end speech recognition in english and mandarin," in *International Conference on Machine Learning*, 2016, pp. 173–182.

[17] Haşim Sak, Andrew Senior, Kanishka Rao, and Françoise Beaufays, "Fast and accurate recurrent neural network acoustic models for speech recognition," in *Proceedings of INTERSPEECH*, 2015.

[18] Yajie Miao, Mohammad Gowayyed, and Florian Metze, "Eesen: End-to-end speech recognition using deep rnn models and wfst-based decoding," in *Automatic Speech Recognition and Understanding (ASRU), 2015 IEEE Workshop on*. IEEE, 2015, pp. 167–174.

[19] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al., "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.

[20] Andrew Maas, Ziang Xie, Dan Jurafsky, and Andrew Ng, "Lexicon-free conversational speech recognition with neural networks," in *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2015, pp. 345–354.

[21] Jan Chorowski, Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, "End-to-end continuous speech recognition using attention-based recurrent nn: first results," *arXiv preprint arXiv:1412.1602*, 2014.

[22] Liang Lu, Xingxing Zhang, Kyunghyun Cho, and Steve Renals, "A study of the recurrent neural network encoder-decoder for large vocabulary speech recognition.," in *INTERSPEECH*, 2015, pp. 3249–3253.

[23] Dzmitry Bahdanau, Jan Chorowski, Dmitriy Serdyuk, Philemon Brakel, and Yoshua Bengio, "End-to-end attention-based large vocabulary speech recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 4945–4949.

[24] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 4960–4964.

[25] K. Audhkhasi, B. Ramabhadran, G. Saon, M. Picheny, and D. Nahamoo, "Direct acoustics-to-word models for English conversational speech recognition," in *Proc. Interspeech*, 2017, pp. 959–963.

[26] H. Soltau, G. Saon, and B. Kingsbury, "The ibm attila speech recognition toolkit," in *Spoken Language Technology Workshop (SLT)*, 2010, pp. 97–102.

[27] Diederik Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[28] H. Bourlard and N. Morgan, "Generalization and parameter estimation in feedforward nets: Some experiments," in *Advances in Neural Information Processing Systems*, 1990, vol. II, pp. 630–637.

[29] Hal Daumé III, "Frustratingly easy domain adaptation," *arXiv preprint arXiv:0907.1815*, 2009.