

# THE MICROSOFT 2017 CONVERSATIONAL SPEECH RECOGNITION SYSTEM

*W. Xiong, L. Wu, F. Alleva, J. Droppo, X. Huang, A. Stolcke*

Microsoft AI and Research

## ABSTRACT

We describe the latest version of Microsoft’s conversational speech recognition system for the Switchboard and CallHome domains. The system adds a CNN-BLSTM acoustic model to the set of model architectures we combined previously, and includes character-based and dialog session aware LSTM language models in rescoring. For system combination we adopt a two-stage approach, whereby acoustic model posteriors are first combined at the senone/frame level, followed by a word-level voting via confusion networks. We also added another language model rescoring step following the confusion network combination. The resulting system yields a 5.1% word error rate on the NIST 2000 Switchboard test set, and 9.8% on the CallHome subset.

**Index Terms**— Conversational speech recognition, CNN, LACE, BLSTM, LSTM-LM, system combination, human parity.

## 1. INTRODUCTION

We have witnessed steady progress in the improvement of automatic speech recognition (ASR) systems for conversational speech, a genre that was once considered among the hardest in the speech recognition community due to its unconstrained nature and intrinsic variability [1]. The combination of deep networks and efficient training methods with older neural modeling concepts [2, 3, 4, 5, 6, 7, 8] have produced steady advances in both acoustic modeling [9, 10, 11, 12, 13, 14, 15] and language modeling [16, 17, 18]. These systems typically use deep convolutional neural network (CNN) architectures in acoustic modeling, and multi-layered recurrent networks with gated memory (long-short-term memory, LSTM [8]) models for both acoustic and language modeling, driving the word error rate on the benchmark Switchboard corpus [19] down from its mid-2000s plateau of around 15% to well below 10%. We can attribute this progress to the neural models’ ability to learn regularities over a wide acoustic context in both time and frequency dimensions, and, in the case of language models, to condition on unlimited histories and learn representations of functional word similarity [20, 21].

Given these developments, we previously carried out an experiment to measure the accuracy of a state-of-the-art conversational speech recognition system against that of professional transcribers. We were trying to answer the question whether machines had effectively caught up with humans in this, originally very challenging, speech recognition task. To measure human error on this task, we submitted the Switchboard evaluation data to our standard conversational speech transcription vendor pipeline (who was left blind to the experiment), postprocessed the output to remove text normalization discrepancies, and then applied the NIST scoring protocol. The resulting human word error was 5.9%, not statistically different from the 5.8% error rate achieved by our ASR system [22]. In a follow-up study [23], we found that qualitatively, too, the human and machine transcriptions were remarkably similar: the same short

function words account for most of the errors, the same speakers tend to be easy or hard to transcribe, and it is difficult for human subjects to tell whether an errorful transcript was produced by a human or ASR. Meanwhile, another research group carried out their own measurement of human transcription error [24], while multiple groups reported further improvements in ASR performance [24, 25]. The IBM/Appen human transcription study employed a more involved transcription process with more listening passes, a pool of transcribers, and access to the conversational context of each utterance, yielding a human error rate of 5.1%. Together with a prior study by LDC [26], we can conclude that human performance, unsurprisingly, falls within a range depending on the level of effort expended.

In this paper we describe a new iteration in the development of our system. Improvements to the acoustic model, the language modeling, rescoring and system combination process combine to yield error rates below the human levels previously measured by us. The remainder of the paper details the components of the system and reports results that elucidate their contributions to the overall performance.

## 2. ACOUSTIC MODELS

### 2.1. Convolutional neural nets

We used two types of CNN model architectures: ResNet and LACE. The residual-network (ResNet) architecture [27] is a standard CNN with added highway connections [28], i.e., a linear transform of each layer’s input to the layer’s output [28, 29]. We apply batch normalization [30] before computing rectified linear unit (ReLU) activations.

The LACE (layer-wise context expansion with attention) model is another CNN architecture first described in [31]. It shares many of the features of time-delay neural networks (TDNNs) [4] but adds attention masking and ResNet-like linear pass-through connections.

### 2.2. Bidirectional LSTM

For our LSTM-based acoustic models we use a bidirectional architecture (BLSTM) [32] without frame-skipping [11]. The core model structure is the LSTM defined in [10]. We found that using networks with more than six layers did not improve the word error rate on the development set, and chose 512 hidden units, per direction, per layer; this gave a reasonable trade-off between training time and final model accuracy.

BLSTM performance was significantly enhanced using a spatial smoothing technique, first described in [22]. Briefly, a two-dimensional topology is imposed on each layer, and activation patterns in which neighboring units are correlated are rewarded.

### 2.3. CNN-BLSTM

A new addition to our system this year is a CNN-BLSTM model inspired by [33]. Unlike the original BLSTM model, we included the context of each time point as an input feature in the model. The context windows was  $[-3, 3]$ , so the input feature has size  $40 \times 7 \times t$ , with zero-padding in the frequency dimension, but not in the time dimension. We first apply three convolutional layers on the features at time  $t$ , and then apply six BLSTM layers to the resulting time sequence, similar to structure of our pure BLSTM model.

### 2.4. Senone set diversity

One standard element of state-of-the-art ASR systems is the combination of multiple acoustic models. Assuming these models are *diverse*, i.e., make errors that are not perfectly correlated, an averaging or voting combination of these models should reduce error. In the past we have relied mainly on different model architectures to produce diverse acoustic models. However, results in [22] for multiple BLSTM models showed that diversity can also be achieved using different sets of senones (clustered subphonetic units). Therefore, we have now adopted a variety of senone sets for all model architectures. Senone sets differ by clustering detail (9k versus 27k senones), as well as two slightly different phone sets and corresponding dictionaries. The standard version is based on the CMU dictionary and phone set (without stress, but including a schwa phone). An alternate dictionary adds specialized vowel and nasal phones used exclusively for filled pauses and backchannel words, inspired by [34]. Combined with set sizes, this gives us a total of four distinct senone sets.

### 2.5. Speaker adaptation

Speaker adaptive modeling in our system is based on conditioning the network on an i-vector [35] characterization of each speaker [36, 37]. A 100-dimensional i-vector is generated for each conversation side (channel A or B of the audio file, i.e., all the speech coming from the same speaker). For the BLSTM systems, the conversation-side i-vector  $v_s$  is appended to each frame of input. For convolutional networks, this approach is inappropriate because we do not expect to see spatially contiguous patterns in the input. Instead, for the CNNs, we add a learnable weight matrix  $W^l$  to each layer, and add  $W^l v_s$  to the activation of the layer before the nonlinearity. Thus, in the CNN, the i-vector essentially serves as a speaker-dependent bias to each layer. For results showing the effectiveness of i-vector adaptation on our models, see [38].

### 2.6. Sequence training

All our models are sequence-trained using maximum mutual information (MMI) as the discriminative objective function. Based on the approaches of [39] and [40], the denominator graph is a full trigram LM over phones and senones. The forward-backward computations are cast as matrix operations, and can therefore be carried out efficiently on GPUs without requiring a lattice approximation of the search space. For details of our implementation and empirical evaluation relative to cross-entropy trained models, see [38].

### 2.7. Frame-level model combination

In our new system we added frame-level combination of senone posteriors from multiple acoustic models. Such a combination of neural

**Table 1.** Acoustic model performance by senone set, model architecture, and for frame-level combinations, using an N-gram LM. The “puhpum” senone sets use an alternate dictionary with special phones for filled pauses.

| Senone set | Architecture | WER devset | WER test |      |
|------------|--------------|------------|----------|------|
|            |              |            | SWB      | CH   |
| 9k         | BLSTM        | 11.5       | 8.3      | 14.3 |
|            | ResNet       | 11.7       | 8.5      | 14.5 |
|            | LACE         | 11.3       | 8.5      | 14.7 |
|            | CNN-BLSTM    | 11.2       | 8.5      | 14.2 |
|            | Combined     | 9.6        | 7.2      | 12.4 |
| 9k puhpum  | BLSTM        | 11.3       | 8.2      | 14.4 |
|            | ResNet       | 11.2       | 8.4      | 14.6 |
|            | LACE         | 11.1       | 8.3      | 14.9 |
|            | CNN-BLSTM    | 11.5       | 8.1      | 14.9 |
|            | Combined     | 9.6        | 7.2      | 12.7 |
| 27k        | BLSTM        | 11.4       | 7.9      | 14.3 |
|            | ResNet       | 11.3       | 8.4      | 14.4 |
|            | LACE         | 11.3       | 8.7      | 14.3 |
|            | CNN-BLSTM    | 11.8       | 8.5      | 14.5 |
|            | Combined     | 9.7        | 7.4      | 12.3 |
| 27k puhpum | BLSTM        | 11.3       | 8.0      | 15.3 |
|            | ResNet       | 11.2       | 8.1      | 14.6 |
|            | LACE         | 11.2       | 8.5      | 14.2 |
|            | CNN-BLSTM    | 11.4       | 8.4      | 14.4 |
|            | Combined     | 9.6        | 7.2      | 12.5 |

acoustic models is effectively just another, albeit more complex, neural model. Frame-level model combination is constrained by the fact that the underlying senone sets must be identical.

Table 1 shows the error rates achieved by various senone set, model architectures, and frame-level combination of all four architectures. The results are based on N-gram language models, and all combinations are equal-weighted.

## 3. LANGUAGE MODELS

### 3.1. Vocabulary size

In the past we had used a relatively small vocabulary of 30,500 words drawn only from in-domain (Switchboard and Fisher corpus) training data. While this yields an out-of-vocabulary (OOV) rate well below 1%, our error rates have reached levels where even small absolute reductions in OOVs could potentially have a significant impact on overall accuracy. We supplemented the in-domain vocabulary with the most frequent words in the out-of-domain sources also used for language model training: the LDC Broadcast News corpus and the UW Conversational Web corpus. Boosting the vocabulary size to 165k reduced the OOV rate (excluding word fragments) on the eval2002 devset from 0.29% to 0.06%. Devset error rate (using the 9k-senones BLSTM+ResNet+LACE acoustic models, see Table 1) dropped from 9.90% to 9.78%.

### 3.2. LSTM-LM rescoring

For each acoustic model our system decodes with a slightly pruned 4-gram LM and generates lattices. These are then rescored with the full 4-gram LM to generate 500-best lists. The N-best lists in turn are then rescored with LSTM-LMs.

Following promising results by other researchers [41, 18], we had already adopted LSTM-LMs in our previous system, with a few enhancements [22]:

- Interpolation of models based on one-hot word encodings (with embedding layer) and another model using letter-trigram word encoding (without extra embedding layer).
- Log-linear combination of forward- and backward-running models.
- Pretraining on the large out-of-domain UW Web corpus (without learning rate adjustment), followed by final training on in-domain data only, with learning rate adjustment schedule.
- Improved convergence through a variation of self-stabilization [42], in which each output vector  $x$  of non-linearities are scaled by  $\frac{1}{4} \ln(1 + e^{4\beta})$ , where a  $\beta$  is a scalar that is learned for each output. This has a similar effect as the scale of the well-known batch normalization technique [30], but can be used in recurrent loops.
- Data-driven learning of the penalty to assign to words that occur in the decoder LM but not in the LSTM-LM vocabulary. The latter consists of all words occurring twice or more in the in-domain data (38k words).

Also, for word-encoded LSTM-LMs, we use the approach from [43] to tie the input embedding and output embedding together.

In our updated system, we add the following additional utterance-scoped LSTM-LM variants:

- A character-based LSTM-LM
- A letter-trigram word-based LSTM-LM using a variant version of text normalization
- A letter-trigram word-based LSTM-LM using a subset of the full in-domain training corpus (a result of holding out a portion of training data for perplexity tuning)

All LSTM-LMs with word-level input use three 1000-dimensional hidden layers. The word embedding layer for the word-based is also of size 1000, and the letter-trigram encoding has size 7190 (the number of unique trigrams). The character-level LSTM-LM uses two 1000-dimensional hidden layers, on top of a 300-dimensional embedding layer.

As before, we build forward and backward running versions of these models, and combine them additively in the log-probability space, using equal weights. Unlike before, we combine the different LSTM-architectures via log-linear combination in the rescoring stage, rather than via linear interpolation at the word level. The new approach is more convenient when the relative weighting of a large number of models needs to be optimized, and the optimization happens jointly with the other knowledge sources, such as the acoustic and pronunciation model scores.

We added one more type of LSTM-LM that represents a more fundamental departure. This LM models the entire dialog sessions instead of individual utterances, but using the entire history of words from the start of the session as conditioning information, along with information about speaker changes and turn overlap. The goal of this session-based LSTM-LM is to capture global coherence in topic and style (entrainment), as well as local cross-turn phenomena such as adjacency pairs; it is described in a companion paper [44].

Table 2 shows perplexities of the various LSTM language models on dev and test sets. The forward and backward versions have very similar perplexities, justifying tying their weights in the eventual score weighting. There are differences between the various input encodings for the utterance-based models, but they are small, on the order of 2-4% relative.

**Table 2.** Perplexities of LSTM-LMs on Switchboard data

| Model structure              | Direction | PPL devset | PPL test |
|------------------------------|-----------|------------|----------|
| Word input, one-hot          | forward   | 50.95      | 44.69    |
|                              | backward  | 51.08      | 44.72    |
| Character input, one-hot     | forward   | 51.66      | 44.24    |
|                              | backward  | 51.92      | 45.00    |
| Word input, letter-trigram   | forward   | 50.76      | 44.55    |
|                              | backward  | 50.99      | 44.76    |
| + alternate text norm        | forward   | 52.08      | 43.87    |
|                              | backward  | 52.02      | 44.23    |
| + alternate training set     | forward   | 50.93      | 43.96    |
|                              | backward  | 50.72      | 44.36    |
| + session-level conditioning | forward   | 37.86      | 35.02    |

Also shown is the effect of session-level modeling, which gives a large perplexity reduction of over 20% over a corresponding utterance-based letter-trigram-encoded LM. For inclusion in the overall system, we built letter-trigram and word-based versions of the session-based LSTM (in both directions). All LSTM-LMs are combined log-linearly at the utterance level (after combining forward and backward variants with equal weights).

## 4. EXPERIMENTAL SETUP

### 4.1. Data

The data sets used for system training are unchanged [22]; they consist of the public and shared data sets used in the DARPA research community. Acoustic training used the English CTS (Switchboard and Fisher) corpora, totalling about 2000 hours. Unlike in previous systems we have reported, we also added CallHome English acoustic data (25 hours after segmentation and alignment, weighted 10-fold), but only in the sequence-training step (to avoid a complete retraining). This addition improved individual acoustic model WER by about 0.5% absolute on CallHome test data.

Language model training used transcripts of the same three CTS corpora, BBN Switchboard-2 transcripts, the LDC Hub4 (Broadcast News) corpus, and the UW conversational web corpus [45]. The Switchboard-1 and Switchboard-2 portions of the NIST 2002 CTS test set were used for tuning and development. Evaluation is carried out on the NIST 2000 CTS test set, with Switchboard (SWB) and CallHome (CH) subsets reported.

### 4.2. Model training

All neural networks in the final system were trained with the Microsoft Cognitive Toolkit, or CNTK [46] on a Linux-based multi-GPU server farm. CNTK allows for flexible model definition, while at the same time scaling efficiently across multiple GPUs and multiple servers. Training times become feasible by parallelizing the stochastic gradient descent (SGD) training with a *1-bit SGD* parallelization technique [47].

We use the CNTK “FsAdaGrad” learning algorithm, which is an implementation of Adam [48]. A typical learning rate is  $3 \times 10^{-6}$ , and learning rates are automatically adjusted with a decrease factor of 0.7. Momentum is set at a constant value of 2500 throughout model training. For individual acoustic models, we find that training converges after 1.5 to 2 passes over the 2000-hour training set. We do not use dropout or gradient noise in our model training, only the aforementioned spatial smoothing technique for BLSTM model training.

**Table 3.** Results for LSTM-LM rescoring on systems selected for combination, the combined system, and confusion network rescoring

| Senone set | Model/combination step        | ngram-LM   |     |      | LSTM-LMs   |     |      |
|------------|-------------------------------|------------|-----|------|------------|-----|------|
|            |                               | WER devset | SWB | CH   | WER devset | SWB | CH   |
| 9k         | BLSTM                         | 11.5       | 8.3 | 14.6 | 9.2        | 6.4 | 12.1 |
| 9k-puhpum  | BLSTM                         | 11.3       | 8.2 | 14.4 | 9.1        | 6.3 | 12.1 |
| 27k        | BLSTM                         | 11.5       | 7.9 | 14.3 | 9.3        | 6.3 | 12.0 |
| 27k-puhpum | BLSTM                         | 11.3       | 8.0 | 15.3 | 9.2        | 6.3 | 12.8 |
| 9k         | BLSTM+ResNet+LACE+CNN-BLSTM   | 9.6        | 7.2 | 12.4 | 7.8        | 5.4 | 10.2 |
| 9k-puhpum  | BLSTM+ResNet+LACE+CNN-BLSTM   | 9.6        | 7.2 | 12.7 | 7.7        | 5.4 | 10.2 |
| 27k        | BLSTM+ResNet+LACE+CNN-BLSTM   | 9.7        | 7.4 | 12.3 | 7.7        | 5.6 | 10.2 |
| 27k-puhpum | BLSTM+ResNet+LACE+CNN-BLSTM   | 9.6        | 7.2 | 12.5 | 7.7        | 5.5 | 10.3 |
| -          | Confusion network combination |            |     |      | 7.3        | 5.2 | 9.8  |
| -          | + LSTM rescoring              |            |     |      | 7.2        | 5.2 | 9.8  |
| -          | + ngram rescoring             |            |     |      | 7.2        | 5.1 | 9.8  |

## 5. SYSTEM COMBINATION AND RESULTS

### 5.1. Confusion network combination

After rescoring all system outputs with all language models, we combine all scores log-linearly and normalize to estimate utterance-level posterior probabilities. All N-best outputs for the same utterance are then concatenated and merged into a single word confusion network (CN), using the SRILM nbest-rover tool [49, 34].

Unlike in our previous system [22], we do not apply estimated, system-level weights to the posterior probabilities estimated from the N-best hypotheses. All systems have equal weight upon combination. The prior work had also shown that a CN combination of all BLSTM system variants (with different senone sets) was highly effective by itself. Consequently, in the present system we combine all four BLSTM systems, as well as the four frame-combined systems.

### 5.2. Confusion network rescoring

As a final processing step, we generate new N-best lists from the confusion networks resulting from system combination. Following [50], these are once more rescored using the N-gram LM, as well as the but also with a subset of the utterance-level LSTM-LMs, and one additional knowledge source. The word log posteriors from the confusion network take the place of the acoustic model scores in this final rescoring step.

Table 3 compares the individual systems that go into the system combination step, before and after rescoring with LSTM-LMs, and then shows the progression of results in the final processing stages, starting with the LM-rescored individual systems, the system combination, and the CN rescoring. The collection of LSTM-LMs (including the session-based LMs) gives very consistent relative error reductions for the individual and frame-combined systems compared to the N-gram LM (about 24% for SWB and 17% for CH). The system combination reduces error by 4% relative over the best individual systems for both SWB and CH. Confusion network rescoring gives a small (0.1% absolute) gain on SWB, but not on CH; this is possibly due to a lack of matched devset data for tuning rescoring weights.

## 6. CONCLUSIONS

We have described the latest iteration of our conversational speech recognition system. The acoustic model was enhanced by adding a CNN-BLSTM system, and the more systematic use of a variety of senone sets, to benefit later system combination. We also switched

to combining different model architectures first at the senone/frame level, resulting in several acoustic combined systems that are then fed into the word-level combination, based on confusion networks. The language model was updated with larger vocabulary (lowering the OOV rate by about 0.2% absolute), additional LSTM-LM variants for rescoring, and most importantly, session-level LSTM-LM that can model global and local coherence between utterances, as well as dialog phenomena. Finally, we added an extra rescoring step where N-best hypotheses generated from the combined confusion network are reweighted with multiple language models, giving a small additional gain for the Switchboard test set. Overall, we have reduced error rate for the Switchboard tasks by 12% relative, from 5.8% for the 2016 system, to now 5.1% for the Switchboard test data, and by 11% relative, to now 9.8% for CallHome test data. We note that these error rates are now below those measured by us previously for a two-pass human transcription pipeline [23].

**Acknowledgments.** We wish to thank our colleagues H. Erdogan, X. He, J. Li, F. Seide, M. Seltzer, and T. Yoshioka for their valued input during system development, and ICSI for assistance with CTS data sets.

## 7. REFERENCES

- [1] S. Greenberg, J. Hollenback, and D. Ellis, "Insights into spoken language gleaned from phonetic transcription of the Switchboard corpus", in *Proc. ICSLP*, 1996.
- [2] F. J. Pineda, "Generalization of back-propagation to recurrent neural networks", *Physical Review Letters*, vol. 59, pp. 2229, 1987.
- [3] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent neural networks", *Neural Computation*, vol. 1, pp. 270–280, 1989.
- [4] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. J. Lang, "Phoneme recognition using time-delay neural networks", *IEEE Trans. Acoustics, Speech, and Signal Processing*, vol. 37, pp. 328–339, 1989.
- [5] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series", *The handbook of brain theory and neural networks*, vol. 3361, pp. 1995, 1995.
- [6] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition", *Neural computation*, vol. 1, pp. 541–551, 1989.
- [7] T. Robinson and F. Fallside, "A recurrent error propagation network speech recognition system", *Computer Speech & Language*, vol. 5, pp. 259–274, 1991.
- [8] S. Hochreiter and J. Schmidhuber, "Long short-term memory", *Neural Computation*, vol. 9, pp. 1735–1780, 1997.

- [9] F. Seide, G. Li, and D. Yu, “Conversational speech transcription using context-dependent deep neural networks”, in *Proc. Interspeech*, pp. 437–440, 2011.
- [10] H. Sak, A. W. Senior, and F. Beaufays, “Long short-term memory recurrent neural network architectures for large scale acoustic modeling”, in *Proc. Interspeech*, pp. 338–342, 2014.
- [11] H. Sak, A. Senior, K. Rao, and F. Beaufays, “Fast and accurate recurrent neural network acoustic models for speech recognition”, in *Proc. Interspeech*, pp. 1468–1472, 2015.
- [12] G. Saon, H.-K. J. Kuo, S. Rennie, and M. Picheny, “The IBM 2015 English conversational telephone speech recognition system”, in *Proc. Interspeech*, pp. 3140–3144, 2015.
- [13] T. Sercu, C. Puhersch, B. Kingsbury, and Y. LeCun, “Very deep multilingual convolutional neural networks for LVCSR”, in *Proc. IEEE ICASSP*, pp. 4955–4959. IEEE, 2016.
- [14] M. Bi, Y. Qian, and K. Yu, “Very deep convolutional neural networks for LVCSR”, in *Proc. Interspeech*, pp. 3259–3263, 2015.
- [15] Y. Qian, M. Bi, T. Tan, and K. Yu, “Very deep convolutional neural networks for noise robust speech recognition”, *IEEE/ACM Trans. Audio, Speech, and Language Processing*, vol. 24, pp. 2263–2276, Aug. 2016.
- [16] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, “Recurrent neural network based language model”, in *Proc. Interspeech*, pp. 1045–1048, 2010.
- [17] M. Sundermeyer, R. Schlüter, and H. Ney, “LSTM neural networks for language modeling”, in *Proc. Interspeech*, pp. 194–197, 2012.
- [18] I. Medennikov, A. Prudnikov, and A. Zatornitskiy, “Improving English conversational telephone speech recognition”, in *Proc. Interspeech*, pp. 2–6, 2016.
- [19] J. J. Godfrey, E. C. Holliman, and J. McDaniel, “Switchboard: Telephone speech corpus for research and development”, in *Proc. IEEE ICASSP*, vol. 1, pp. 517–520. IEEE, 1992.
- [20] Y. Bengio, H. Schwenk, J.-S. Senécal, F. Morin, and J.-L. Gauvain, “Neural probabilistic language models”, in *Studies in Fuzziness and Soft Computing*, vol. 194, pp. 137–186. 2006.
- [21] T. Mikolov, W.-t. Yih, and G. Zweig, “Linguistic regularities in continuous space word representations”, in *HLT-NAACL*, vol. 13, pp. 746–751, 2013.
- [22] W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, and G. Zweig, “Achieving human parity in conversational speech recognition”, Technical Report MSR-TR-2016-71, Microsoft Research, Oct. 2016, <https://arxiv.org/abs/1610.05256>.
- [23] A. Stolcke and J. Droppo, “Comparing human and machine errors in conversational speech transcription”, in *Proc. Interspeech*, pp. 137–141, Stockholm, Aug. 2017.
- [24] G. Saon, G. Kurata, T. Sercu, K. Audhkhasi, S. Thomas, D. Dimitriadis, X. Cui, B. Ramabhadran, M. Picheny, L.-L. Lim, B. Roomi, and P. Hall, “English conversational telephone speech recognition by humans and machines”, in *Proc. Interspeech*, pp. 132–136, Stockholm, Aug. 2017.
- [25] K. J. Han, S. Hahn, B.-H. Kim, J. Kim, and I. Lane, “Deep learning-based telephony speech recognition in the wild”, in *Proc. Interspeech*, pp. 1323–1327, Stockholm, Aug. 2017.
- [26] M. L. Glenn, S. Strassel, H. Lee, K. Maeda, R. Zakhary, and X. Li, “Transcription methods for consistency, volume and efficiency”, in *Proc. 7th Intl. Conf. on Language Resources and Evaluation*, pp. 2915–2920, Valletta, Malta, 2010.
- [27] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition”, arXiv preprint arXiv:1512.03385, 2015.
- [28] R. K. Srivastava, K. Greff, and J. Schmidhuber, “Highway networks”, *CoRR*, vol. abs/1505.00387, 2015.
- [29] P. Ghahremani, J. Droppo, and M. L. Seltzer, “Linearly augmented deep neural network”, in *Proc. IEEE ICASSP*, pp. 5085–5089. IEEE, 2016.
- [30] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift”, *Proceedings of Machine Learning Research*, vol. 37, pp. 448–456, 2015.
- [31] D. Yu, W. Xiong, J. Droppo, A. Stolcke, G. Ye, J. Li, and G. Zweig, “Deep convolutional neural networks with layer-wise context expansion and attention”, in *Proc. Interspeech*, pp. 17–21, 2016.
- [32] A. Graves and J. Schmidhuber, “Framewise phoneme classification with bidirectional LSTM and other neural network architectures”, *Neural Networks*, vol. 18, pp. 602–610, 2005.
- [33] T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, “Convolutional, long short-term memory, fully connected deep neural networks”, in *Proc. IEEE ICASSP*, 2015.
- [34] A. Stolcke et al., “The SRI March 2000 Hub-5 conversational speech transcription system”, in *Proceedings NIST Speech Transcription Workshop*, College Park, MD, May 2000.
- [35] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, “Front-end factor analysis for speaker verification”, *IEEE Trans. Audio, Speech, and Language Processing*, vol. 19, pp. 788–798, 2011.
- [36] G. Saon, H. Soltau, D. Nahamoo, and M. Picheny, “Speaker adaptation of neural network acoustic models using i-vectors”, in *IEEE Speech Recognition and Understanding Workshop*, pp. 55–59, 2013.
- [37] G. Saon, T. Sercu, S. J. Rennie, and H. J. Kuo, “The IBM 2016 English conversational telephone speech recognition system”, in *Proc. Interspeech*, pp. 7–11, Sep. 2016.
- [38] W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, and G. Zweig, “The Microsoft 2016 conversational speech recognition system”, in *Proc. IEEE ICASSP*, pp. 5255–5259, 2017.
- [39] S. F. Chen, B. Kingsbury, L. Mangu, D. Povey, G. Saon, H. Soltau, and G. Zweig, “Advances in speech transcription at IBM under the DARPA EARS program”, *IEEE Trans. Audio, Speech, and Language Processing*, vol. 14, pp. 1596–1608, 2006.
- [40] D. Povey, V. Peddinti, D. Galvez, P. Ghahremani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, “Purely sequence-trained neural networks for ASR based on lattice-free MMI”, in *Proc. Interspeech*, pp. 2751–2755, 2016.
- [41] M. Sundermeyer, H. Ney, and R. Schlüter, “From feedforward to recurrent LSTM neural networks for language modeling”, *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, pp. 517–529, Mar. 2015.
- [42] P. Ghahremani and J. Droppo, “Self-stabilized deep neural network”, in *Proc. IEEE ICASSP*, pp. 5450–5454. IEEE, 2016.
- [43] O. Press and L. Wolf, “Using the output embedding to improve language models”, arXiv preprint arXiv:1608.05859, 2016.
- [44] W. Xiong, L. Wu, and A. Stolcke, “Conversational session-level language modeling with LSTMs”, Submitted to IEEE ICASSP-2018, 2017.
- [45] I. Bulyko, M. Ostendorf, and A. Stolcke, “Getting more mileage from web text sources for conversational speech language modeling using class-dependent mixtures”, in M. Hearst and M. Ostendorf, editors, *Proceedings of HLT-NAACL 2003, Conference of the North American Chapter of the Association of Computational Linguistics*, vol. 2, pp. 7–9, Edmonton, Alberta, Canada, Mar. 2003. Association for Computational Linguistics.
- [46] Microsoft Research, “The Microsoft Cognition Toolkit (CNTK)”, <https://cntk.ai>.
- [47] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, “1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs”, in *Proc. Interspeech*, pp. 1058–1062, 2014.
- [48] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization”, *Proceedings 3rd International Conference for Learning Representations*, arXiv preprint arXiv:1703.02136, May 2015.
- [49] A. Stolcke, “SRILM—an extensible language modeling toolkit”, in *Proc. Interspeech*, pp. 2002–2005, 2002.
- [50] S. Bangalore, G. Bordel, and G. Riccardi, “Computing consensus translation from multiple machine translation systems”, in *Proceedings IEEE Workshop on Automatic Speech Recognition and Understanding*, pp. 351–354, Madonna di Campiglio, Italy, Dec. 2001.