
Training Neural Speech Recognition Systems with Synthetic Speech Augmentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Building an accurate automatic speech recognition (ASR) system requires a large
2 dataset that contains many hours of labeled speech samples produced by a diverse
3 set of speakers. The lack of such open free datasets is one of the main issues
4 preventing advancements in ASR research. To address this problem, we propose to
5 augment a natural speech dataset with synthetic speech. We train very large end-to-
6 end neural speech recognition models using the LibriSpeech dataset augmented
7 with synthetic speech. These new models achieve state of the art Word Error Rate
8 (WER) for character-level based models without an external language model.

9 1 Introduction

10 There has been a large amount of success in using neural networks (NN) for automatic speech
11 recognition (ASR). Classical ASR systems use complicated pipelines with many heavily engineered
12 processing stages, including specialized input features, acoustic models, and Hidden Markov Models
13 (HMMs). Deep NNs have traditionally been used for acoustic modeling (Weibel et al., 1989; Hinton
14 et al., 2012). A key breakthrough occurred when a state of the art ASR system, Deep Speech, was
15 built using an end-to-end deep learning approach (Hannun et al., 2014). This model replaced both
16 acoustic modeling, and HMMs with very deep neural networks and was able to directly translate
17 spectrograms into English text transcriptions. This end-to-end approach was extended in follow-up
18 papers (Amodei et al., 2016; Collobert et al., 2016). Recent advances in ASR have further improved
19 upon these models using even more advanced techniques such as replacing n-gram language models
20 with a neural language model, usually in the form of a recurrent neural network (RNN) (Zeyer et al.,
21 2018; Povey et al., 2018; Han et al., 2018).

22 As opposed to making neural networks more complex, we were interested in an orthogonal direction
23 of study: whether we can improve quality by creating larger models. In order to train such large
24 models, deep NNs require a vast quantity of data to be available. However the lack of a large public
25 speech dataset blocked us from successfully building large NN models for ASR. We were inspired
26 by recent work in improving translation systems using synthetic data (Sennrich et al., 2015), and as
27 such, we decided to augment speech with synthetic data.¹ Furthermore, given the recent impressive
28 improvement in neural speech synthesis models (van den Oord et al., 2016; Shen et al., 2018), it
29 becomes cheap to generate high quality speech with varying prosody.

30 We show that by using synthetic speech created from a neural speech synthesis model, we can improve
31 ASR performance compared to models trained using only LibriSpeech data (Panayotov et al., 2015).
32 By naively increasing the depth of the model, we observe that the synthetic data allows us to achieve
33 state of the art WER using a greedy decoder.

¹ Synthesized speech has previously been used to improve speech recognition for low-resource languages (Rygaard, 2015).

2 Synthetic Speech Dataset

We use the Tacotron-2 like model from the OpenSeq2Seq² toolkit (Kuchaiev et al., 2018) and add Global Style Tokens (GST) (Wang et al., 2018) to learn multiple speaker identities. Tacotron-2 with GST (T2-GST) was trained on the MAILABS English-US dataset (M-AILABS, 2018) with approximately 100 hours of audio recorded by 3 speakers. T2-GST was able to learn all 3 different speaking styles and different accents that they use to portray different characters.

Using the T2-GST model, we created a fully synthetic version of the LibriSpeech training audio. In order to produce audio, T2-GST requires a spectrogram used for the style and the audio transcription. For the transcription, we took the transcripts from the train-clean-100, train-clean-360, and train-other-500 LibriSpeech splits and randomly paired them with style spectrograms from the MAILABS English-US dataset. At the end, we had a dataset that was the same size as the original training portion of the LibriSpeech dataset but spoken in the tones of the speakers from the MAILABS dataset.

The audio from the T2-GST model could be further controlled by the amount of dropout in the prenet of the decoder. By decreasing the dropout rate, we found that we could slightly distort the audio. The main difference that we noticed was that the lower the dropout rate, the faster the resulting audio would sound. Using this observation, we further increased the size of the synthetic dataset. Thus, we used 46%, 48%, and 50% for the dropout and created a synthetic speech dataset that was 3 times as large as the LibriSpeech training dataset.

3 Training Speech Recognition with Synthetic Data

3.1 Neural Speech Recognition Models

Our speech recognition model is an end-to-end neural network that takes logarithmic mel-scale spectrograms as inputs and produces characters. We use a deep convolutional NN model, which we will further address as Wave2Letter+ (w2lp)³. It is based on Wav2Letter (Collobert et al., 2016) except:

- We use ReLU instead of Gated Linear Unit
- We use batch normalization instead of weight normalization
- We add residual connections between convolutional blocks
- We use Connectionist Temporal Classification loss instead of Auto Segmentation Criterion
- We use Layer-wise Adaptive Rate clipping (LARC) for gradient clipping

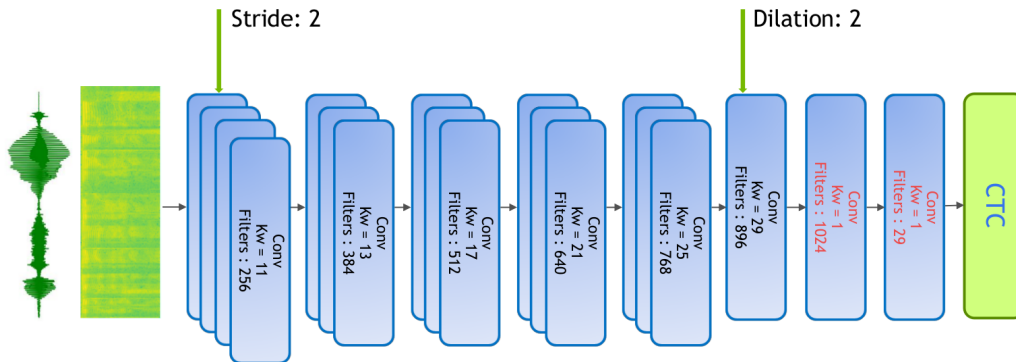


Figure 1: Base Wave2Letter+ model.

²We used OpenSeq2Seq both to create a synthetic speech dataset and build ASR.

³More model details: <https://nvidia.github.io/OpenSeq2Seq/html/speech-recognition/wave2letter.html>

63 The base Wave2Letter+ model has 19 convolutional layers:

- 64 • 1 pre-processing layer at the beginning of the network
- 65 • 15 layers consists of 5 blocks of 3 repeating convolutional layers
- 66 • 3 post-processing layers at the end of the network

67 We made the base model deeper and experimented with 24, 34, 44, and 54 layer networks. The 24
68 layer networks consists of 5 blocks of 4 repeating convolutional layers. The 34, 44, and 54 layer
69 network consists of 10 blocks of 3, 4, and 5 repeating layers respectively.

70 3.2 Word Error Rate Improvement using Synthetic Augmentation

71 The training dataset was created by combining the synthetic data with the original training data from
72 LibriSpeech. Despite the synthetic data being larger than the natural data, we found it most helpful to
73 sample natural and synthetic data at a 50/50 rate.

74 Our best performing model, the 54 layer network, currently has a word error rate of 4.32% on
75 test-clean and a word error rate of 14.08% on test-other by greedily choosing the most probable
76 character at each step without any language model rescoreing. To the best of our knowledge, the
77 previous best performing model without using any language model achieves 4.87% on test-clean
78 and 15.39% on test-other (Zeyer et al., 2018). The complete results can be found in Table 1. Models
79 trained on the combined dataset outperform those trained on original LibriSpeech. The augmented
80 dataset improves results on test-clean by 0.15, and 0.44 for the 24 and 34 layer models. For test-other,
81 we see an improvement of 0.08, and 0.74 for the 24 and 34 layer models.

82 By using beam search with beam width 128 and the 4-gram OpenSLR⁴ language model rescoreing,
83 we improved our WER on test-other to 12.21% on the 54 layer model which is better than previous
84 public 4-gram language models and comparable to LSTM language models (Zeyer et al., 2018).

Model	Dataset Used	Dev		Test	
		Clean	Other	Clean	Other
attention-Zeyer et al.	LibriSpeech	4.87	14.37	4.87	15.39
w2lp-24	LibriSpeech	5.44	16.57	5.31	17.09
w2lp-24	Combined	5.12	16.25	5.16	17.01
w2lp-34	LibriSpeech	5.10	15.49	5.10	16.21
w2lp-34	Combined	4.60	14.98	4.66	15.47
w2lp-44	Combined	4.24	13.87	4.36	14.37
w2lp-54	Combined	4.32	13.74	4.32	14.08

Table 1: Greedy WER on LibriSpeech for Different Models and Datasets

85 3.3 How To Mix Natural and Synthetic Speech

86 We performed a number of additional experiments to find the best sampling ratio between synthetic
87 data and LibriSpeech. We tested training on only LibriSpeech, a 50/50 split, a 33/66 split, and a pure
88 synthetic dataset. All tests were done on the 34 layer model. The results are shown in Table 2.

Model	Sampling Ratio (Natural/Synthetic)	Dev		Test	
		Clean	Other	Clean	Other
w2lp-34	Natural	5.10	15.49	5.10	16.21
w2lp-34	50/50	4.60	14.98	4.66	15.47
w2lp-34	33/66	4.91	15.18	4.81	15.81
w2lp-34	Synthetic	51.39	80.27	49.80	81.78

Table 2: Greedy WER for Different Ratios Between Natural and Synthetic Datasets

⁴The LM can be found here: www.openslr.org/11

89 Despite the larger amount of synthetic data, the synthetic dataset fails to capture the larger variety of
 90 LibriSpeech. We believe that this effect could be moderated if a speech synthesis model with larger
 91 speaker variety was used as opposed to the current 3 speaker speech synthesis model. A 50/50 split
 92 between the natural and synthetic seems to be a good ratio for our dataset.

93 3.4 Traditional Speech Augmentation vs Synthetic Speech

94 Adding synthetic data proved to be better regularization than standard regularization techniques. In
 95 addition to dropout which is employed for all models, OpenSeq2Seq supports speech augmentation
 96 such as adding noise and time stretching. Using these 3 techniques, we tested 4 additional models.
 97 We tested 2 larger dropout factors, and on top on this, we tested whether speech augmentation would
 98 improve performance. Since the dropout factor varies by layer, we multiply the local dropout keep
 99 probabilities by a global dropout keep factor. All tests were done on the 34 layer model. The tests
 100 and results are detailed in Table 3.

101 A slightly larger dropout resulted in minor improvement in WER. The effects of speech augmentation
 102 seem to be negligible or, in the case of large dropout, make WER worse. Adding synthetic data
 103 significantly outperforms all other methods of regularization.

Model	Dropout Keep Factor	Time Stretch Factor	Noise (dB)	Dev		Test	
				Clean	Other	Clean	Other
LibriSpeech	None	None	None	5.10	15.49	5.10	16.21
Dropout	0.9	None	None	5.01	15.15	5.15	15.70
Dropout + Aug	0.9	0.05	[-90, -60]	5.07	15.00	5.02	15.83
Dropout	0.75	None	None	5.46	15.77	5.39	16.62
Dropout + Aug	0.75	0.1	[-90, -60]	5.80	16.33	5.72	17.41
Combined	None	None	None	4.60	14.98	4.66	15.47

Table 3: Greedy WER on Using Different Regularization Techniques

104 4 Conclusion and Future Plans

105 Using synthetic data is an effective way to build large neural speech recognition systems. The
 106 synthetic data should be combined with the natural data in the correct ratio to obtain best results.
 107 With this method, we achieved a WER of 4.32% on test-clean and a WER of 14.08% on test-other
 108 using a greedy decoder. This is the current state of the art on character-level based greedy decoding.
 109 Furthermore, using a language model and a beam search width of 128, we get 12.21% WER on
 110 test-other.

111 For future studies, we are interested in creating a larger synthetic dataset with noise. For now, we have
 112 restricted ourselves to take transcriptions from the training subsets of LibriSpeech, but the speech
 113 synthesis models are general enough to accept any transcript. It would be interesting to scrape text
 114 from another source and add to the training set additional phrases not found in LibriSpeech.

115 Acknowledgments

116 We would like to thank Rafael Valle, Ryan Leary, Igor Gitman, Oleksii Kuchaiev and Ujval Kapasi
 117 for helpful conversation and advice through this project.

118 References

119 Dario Amodei, Rishita Anubhai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Jingdong
 120 Chen, Mike Chrzanowski, Adam Coates, Greg Diamos, Erich Elsen, Jesse Engel, Linxi Fan,
 121 Christopher Fougner, Awni Y. Hannun, Billy Jun, Tony Han, Patrick LeGresley, Xiangang Li,
 122 Libby Lin, Sharan Narang, Andrew Y. Ng, Sherjil Ozair, Ryan Prenger, Sheng Qian, Jonathan
 123 Raiman, Sanjeev Satheesh, David Seetapun, Shubho Sengupta, Chong Wang, Yi Wang, Zhiqian
 124 Wang, Bo Xiao, Yan Xie, Dani Yogatama, Jun Zhan, and Zhenyao Zhu. 2016. Deep speech 2 :
 125 End-to-end speech recognition in english and mandarin. In *ICML*.

- 126 Ronan Collobert, Christian Puhersch, and Gabriel Synnaeve. 2016. Wav2letter: an end-to-end convnet-
127 based speech recognition system. *arXiv preprint arXiv:1609.03193*.
- 128 Kyu J. Han, Akshay Chandrashekar, Jungsuk Kim, and Ian R. Lane. 2018. The CAPIO 2017
129 conversational speech recognition system. *CoRR*, abs/1801.00059.
- 130 Awni Y. Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan
131 Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and Andrew Y. Ng. 2014. Deep
132 speech: Scaling up end-to-end speech recognition. *CoRR*, abs/1412.5567.
- 133 Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew
134 Senior, Vincent Vanhoucke, Patrick Nguyen, Brian Kingsbury, and Tara Sainath. 2012. Deep
135 neural networks for acoustic modeling in speech recognition. *IEEE Signal Processing Magazine*.
- 136 Oleksii Kuchaiev, Boris Ginsburg, Igor Gitman, Vitaly Lavrukhin, Carl Case, and Paulius Micike-
137 vicius. 2018. Openseq2seq: extensible toolkit for distributed and mixed precision training of
138 sequence-to-sequence models. <https://arxiv.org/abs/1805.10387>.
- 139 M-AILABS. 2018. The M-AILABS Speech Dataset. [http://www.m-ailabs.bayern/en/
140 the-mailabs-speech-dataset/](http://www.m-ailabs.bayern/en/the-mailabs-speech-dataset/). Accessed: 2018-10-09.
- 141 Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,
142 Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. Wavenet: A generative
143 model for raw audio. *CoRR*, abs/1609.03499.
- 144 Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an asr
145 corpus based on public domain audio books. In *Acoustics, Speech and Signal Processing (ICASSP),
146 2015 IEEE International Conference on*, pages 5206–5210. IEEE.
- 147 Daniel Povey, Gaofeng Cheng, Yiming Wang, Ke Li, Hainan Xu, Mahsa Yarmohammadi, and Sanjeev
148 Khudanpur. 2018. Semi-orthogonal low-rank matrix factorization for deep neural networks. In
149 *Interspeech*.
- 150 Luise Valentin Rygaard. 2015. Using synthesized speech to improve speech recognition for low-
151 resource languages. <https://parasol.tamu.edu/dreu2015/Rygaard/report.pdf>.
- 152 Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation
153 models with monolingual data. *CoRR*, abs/1511.06709.
- 154 Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng
155 Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. 2018. Natural tts synthesis by conditioning
156 wavenet on mel spectrogram predictions. In *2018 IEEE International Conference on Acoustics,
157 Speech and Signal Processing (ICASSP)*, pages 4779–4783. IEEE.
- 158 Yuxuan Wang, Daisy Stanton, Yu Zhang, R. J. Skerry-Ryan, Eric Battenberg, Joel Shor, Ying Xiao,
159 Fei Ren, Ye Jia, and Rif A. Saurous. 2018. Style tokens: Unsupervised style modeling, control and
160 transfer in end-to-end speech synthesis. *CoRR*, abs/1803.09017.
- 161 Alexander Weibel, Toshiyuki Hanazawa, Geoffrey Hinton, Kiyohiro Shirano, and kevin Lang. 1989. A
162 time-delay neural network architecture for isolated word recognition. *IEEE Trans. on Acoustics,
163 Speech and Signal Processing*.
- 164 Albert Zeyer, Kazuki Irie, Ralf Schlüter, and Hermann Ney. 2018. Improved training of end-to-end
165 attention models for speech recognition. *CoRR*, abs/1805.03294.