

SPEECH SYNTHESIS AS AUGMENTATION FOR LOW-RESOURCE ASR

Deblin Bagchi*, Shannon Wotherspoon*, Zhuolin Jiang, Prasanna Muthukumar†

Raytheon BBN, Cambridge, MA-02138, USA

{dbagchi, swothers, zjiang, pmuthuku}@bbn.com

ABSTRACT

Speech synthesis might hold the key to low-resource speech recognition. Data augmentation techniques have become an essential part of modern speech recognition training. Yet, they are simple, naive, and rarely reflect real-world conditions. Meanwhile, speech synthesis techniques have been rapidly getting closer to the goal of achieving human-like speech. In this paper, we investigate the possibility of using synthesized speech as a form of data augmentation to lower the resources necessary to build a speech recognizer. We experiment with three different kinds of synthesizers: statistical parametric, neural, and adversarial. Our findings are interesting and point to new research directions for the future.

Index Terms— speech synthesis, data augmentation, speech recognition, low resource languages

1. INTRODUCTION

The modern speech processing system has a voracious appetite for data. The general trend in automatic speech recognition (ASR) systems today has been to create algorithms that can make use of increasingly larger datasets. These efforts have undoubtedly yielded impressive results. The ability to learn effectively from large datasets has directly corresponded with recognition accuracy. Far from being unique to speech recognition, ASR's counterpart, speech synthesis has also followed a similar trend. The original Blizzard challenge dataset [1] consisted of a one-hour corpus. The most recent version [2] uses corpora several times larger.

While this trend of ever larger datasets has led to great advances in speech processing abilities, it has also ignored the needs of several end-users. Large sized corpora are simply not available for a large variety of tasks. The most common use case with this constraint is speech recognition when dealing with non-mainstream languages. Igbo, for instance, is spoken by 47 million speakers but no large corpus exists, and creating a several thousand hour Igbo corpus for ASR will be prohibitively expensive. Arguably, older techniques like Hidden Markov Models [3] would still be applicable when dealing with smaller corpora. However, it is important that we do not exclude poorly resourced languages from enjoying the advances of modern speech processing technologies.

In this paper, we will therefore pursue a goal directly *opposite* the general trend of papers today. We will attempt

to reduce the amount of data required to build a viable speech recognizer while still using modern neural network based ASR techniques. We attempt this using speech synthesis techniques [4] as a form of data augmentation. The reasoning behind this decision is that synthesizers typically require orders of magnitude less data than recognizers do. After all, statistical parametric speech synthesizers have even been built on as little as 30 minutes of speech data [5]. Speech synthesizers also offer a means of data augmentation that is less naive compared to standard augmentation techniques. For instance, a common form of data augmentation is to speed up or slow down existing speech. While this technique might improve accuracy, it does not reflect the true way that humans speed up or slow down their speech. A synthesis-based augmentation technique, on the other hand, is more likely to reflect rate changes more accurately. We therefore believe that synthesis-based augmentation should yield even better results.

We will investigate three different forms of speech synthesis and their performance when used as a means of data augmentation for ASR. The first is the classic statistical parametric speech synthesis [4]. The second is the popular neural synthesizer, Tacotron2 [6]. The third is the adversarial synthesizer, WGANSing [7]. We report our experiences with all three synthesizers.

2. PRIOR WORK

Low-resource ASR and TTS have a rich and varied history. Technically, any ASR paper from 10 years ago or earlier can be considered a viable approach for low-resource ASR. However, the closest related work to ours is the paper from Rosenberg et al [8]. The authors describe an approach where the Tacotron speech synthesizer is used to augment real speech with the goal of increasing lexical and acoustic diversity. The major difference between Rosenberg et al. and our work is in the quantity of data used. Keeping in mind our goal of using the least amount of data possible, we use substantially smaller subsets of the LibriSpeech corpus [9]. Because we restricted ourselves to using as little data as possible, we also had to resort to synthesizers beyond Tacotron that Rosenberg et al. use.

Another closely related work is LRSpeech by Xu et al [10]. The authors target low-resource scenarios like we do but use a high-resource language to bootstrap ASR and TTS systems for the low-resource language. We, on the other hand, assume that no corpus is available apart from that of the low-

*Equal contribution

†Corresponding author: Prasanna Muthukumar, pmuthuku@bbn.com

resource language. We see our approach as complementary to that of Xu et al., and hope our approach is an interesting alternative.

3. EXPERIMENTS

3.1. Speech recognizer

Our speech recognition models are trained using the BBN speech recognition system, Sage [11]. Sage is an extension of the popular Kaldi speech recognition toolkit [12]. Our ASR systems are multilingual initialized [13] hybrid TDNN-F models [14], trained for one epoch of lattice-free MMI (LF-MMI) and two epochs of sMBR. For language modeling, we train word-level trigram models. All our experiments in this paper will focus on the acoustic model since text data is typically easier to obtain than labeled audio. We use 100 hours of Librispeech transcripts for language model data in all cases, and vary the amount of audio data as required by the experiment.

3.2. Statistical Parametric synthesis

For our experiments using statistical parametric speech synthesis, we use the Clustergen synthesizer [15]. Clustergen uses the Festival system [16] to convert text into a sequence of phonemes, and uses random forests to predict mel-frequency cepstral co-efficients that correspond to the phonemes.

While Clustergen generally produces good quality speech, we use it differently from its original intent. Clustergen is heavily optimized towards a single speaker corpus. However, most ASR corpora are multi-speaker, and only possess a few minutes of speech for each speaker. It is impossible to build a synthesis model for each speaker with such little data. We therefore attempted to build average models for groups of speakers. We identified similar speakers by clustering i-vectors [17]. We pretended that all the speech from a particular cluster was from the same speaker, and built a Clustergen model for it. These average models did not produce high quality synthesis but resulted in some interesting findings anyway. They pointed out several flaws in Mel Cepstral Distortion (MCD), the objective metric we were using to assess synthesis quality. We discuss these flaws in section 4.

Table 1 shows the results of using Clustergen as a data augmentation system. Here, Clustergen was trained on the specified amount of real speech and then used to synthesize the specified amount of synthetic speech. The ‘MCD < 5’ and ‘MCD < 6’ rows correspond to experiments where we left out synthetic voices that had MCD scores higher than the number.

Table 2 shows the results of an experiment where we trained on either purely real speech or purely synthetic speech. The results in the second row were generated by building a separate Clustergen model for each speaker no matter how little the available data. The third row in the table is the result of clustering similar speakers using i-vectors. The results in both tables indicate that Clustergen’s speech

Table 1. ASR performance when using synthetic speech as additional training data

Training data	WER
20h real	12.7
20h real + 20h synth	12.6
20h real + 20h synth (MCD < 5)	12.7
20h real + 20h synth (MCD < 6)	12.7
20h real + 60h synth	13.0
20h real + 80h synth	12.8
40h real	12.0
40h real + 60h synth	13.0
80h real	11.4
80h real + 20h synth	11.5

Table 2. ASR performance when training on either purely real or purely synthetic speech

Training data	WER
80h real	11.4
80h synth, unclustered, 50 voices	36.2
80h synth, clustered, 5 voices	47.1

quality was not sufficient to be useful as a data augmentation method.

For neural TTS and adversarial TTS, we used Clustergen as a baseline. Since ASR training is computationally expensive, we did not train the ASR system whenever any synthesizer performed worse than the parametric synthesizer.

3.3. Neural TTS

Tacotron2 [6] is a purely neural speech synthesizer that consists of a recurrent sequence-to-sequence neural network that takes as input character embeddings and predicts mel spectrograms. We use Waveglow [18] as a vocoder to convert the mel spectrograms to raw speech. Tacotron2 has been shown to produce excellent quality speech in the past, but our primary reason for choosing this synthesizer and waveglow was because there was open-source code available.

We were successful in producing high quality synthesis from Tacotron2 by training on the 24 hour LJ Speech dataset [19], but had difficulty achieving similar results on any 1 hour CMU Arctic dataset [20]. Despite our best efforts, the best Tacotron2 could do on the one-hour corpus was babble in a voice that strongly resembled the original speaker. This test is important because CMU Arctic is still larger and cleaner than any single voice we had for our ASR datasets. We therefore had to rule out any possibility of training a Tacotron2 model exclusive for each voice in an ASR corpus. (Building an exclusive model for each voice in an ASR corpus would have been terribly expensive computationally anyway).

To overcome this difficulty, we extended Tacotron2 to

support multi-speaker training by incorporating speaker embeddings as input in addition to the character embeddings. We experimented with both one-hot embeddings as well as i-vectors. One-hot embeddings gave us marginally better results than i-vectors but neither could match up to even the quality of Clustergen. Nor was there an easy way of quantifying the difference in quality between the i-vectors and one-hot embeddings. Since the quality of synthesis on our corpus was significantly worse than using classical parametric synthesis, we did not attempt our augmentation experiments.

3.4. Adversarial TTS

In addition to the classical parametric and neural speech synthesizers, we also tested an adversarial synthesizer based on Generative Adversarial Networks (GANs). Adversarial learning is a new machine learning technique that has had tremendous success in generating high-dimensional data. These techniques are still in the nascent stages for speech synthesis, but nevertheless show great promise because of the success they have had in related fields such as image generation.

The toolkit we used for our experiments was actually a GAN-based singing synthesizer called WGANSing [7]. We chose WGANSing because it was the only open-source toolkit we could find capable of open-vocabulary speech synthesis. The architecture of this system is the same as DCGANs [21] and is trained using the Wasserstein GAN algorithm.

WGANSing takes as input a block of frame-wise linguistic features and singer identity features, and outputs vocoder features that correspond to the block. The most problematic for us among the input features were the precise phone durations and the pitch contours required. These features make sense for synthesizing singing, but are very difficult to generate for regular speech. We were therefore forced to modify the toolkit to avoid conditioning on these troublesome features. However, the resulting synthesis quality without these features was poor compared to the other synthesizers discussed in this paper. We suspect that this sub-par performance is caused by the low amounts of training data we use. Low-resources are a requirement for our task. Unfortunately, this requirement might also be the Achilles heel of adversarial ML.

4. DISCUSSION

One major setback we faced in our series of experiments is the failure of Mel Cepstral Distortion (MCD) as a metric. In standard speech synthesis experiments, MCD has issues but is generally reliable. For instance, a rule of thumb is that an MCD reduction of 0.12 is equivalent to doubling the amount of training data [22]. Unfortunately, intuitive rules and beliefs such as these only ended up being true for the standard speech synthesis training scenario of clean, single-speaker data. When training on noisy, multi-speaker data, not

only was MCD a poor reflection of quality, but it was also extremely misleading. In informal tests, we found several instances where the metric indicated high MCD but the synthesized speech was understandable. We also found several examples where the synthesized speech was unintelligible but the measured MCD was low. Clearly, single-speaker objective metrics do not extrapolate well to multi-speaker training scenarios.

Another major pain point for us was the gordian knot that was Tacotron2. We faced tremendous engineering and research challenges in getting sensible performance from this synthesizer. The engineering challenges were primarily related to the amount of compute required (at least the implementation we used required this). Both Tacotron2 and Waveglow could only be run on the fastest GPUs we had available, and then still took between days to weeks to converge. While we managed to overcome these through a lot of work, the research challenges proved more insurmountable. We were able to successfully replicate the default experiments in the open-source Tacotron2 toolkit, but we had little success when running the codebase on our own smaller dataset. The inscrutable nature of modern neural networks also made it difficult to extend the system to support multiple speakers. Despite being experts in speech synthesis and neural networks, progress was extremely difficult. We hope that future work by us and others relieves researchers of this burden, and makes it easier to achieve positive results.

5. CONCLUSION

The results of our experiments have mostly been disappointing. It appears that modern speech synthesis has not yet advanced to be of sufficient use in training low-resource speech recognizers. Nevertheless we still believe that as speech synthesizers increase in quality, this approach will eventually become viable. We also hope that our experiments point to interesting future directions with a higher potential for success. For instance, every synthesizer we explored in our paper has been heavily tuned for the human ear. It is highly likely that new synthesizers will need to be designed and built with the explicit goal of acting as sources of augmented data.

6. ACKNOWLEDGEMENT

We thank Sean Colbath, Ilana Heintz, and Ron Coleman for all their support in this endeavor.

7. REFERENCES

- [1] Alan W Black and Keiichi Tokuda, “The blizzard challenge-2005: Evaluating corpus-based speech synthesis on common datasets,” in *Ninth European Conference on Speech Communication and Technology*, 2005.
- [2] Zhenhua Ling and Simon King, “The blizzard challenge 2020: Evaluating corpus-based speech synthesis on common datasets,” <https://www.synsig.org/images/e/ee/Blizzard2020-CallforParticipation.pdf>.
- [3] Lawrence Rabiner and B Juang, “An introduction to hidden markov models,” *ieee assp magazine*, vol. 3, no. 1, pp. 4–16, 1986.
- [4] Heiga Zen, Keiichi Tokuda, and Alan W Black, “Statistical parametric speech synthesis,” *speech communication*, vol. 51, no. 11, pp. 1039–1064, 2009.
- [5] Alan W Black and Prasanna Kumar Muthukumar, “Random forests for statistical speech synthesis,” in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [6] Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al., “Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4779–4783.
- [7] Pritish Chandna, Merlijn Blaauw, Jordi Bonada, and Emilia Gómez, “WGANSing: A multi-voice singing voice synthesizer based on the wasserstein-gan,” in *2019 27th European Signal Processing Conference (EUSIPCO)*. IEEE, 2019, pp. 1–5.
- [8] Andrew Rosenberg, Yu Zhang, Bhuvana Ramabhadran, Ye Jia, Pedro Moreno, Yonghui Wu, and Zelin Wu, “Speech recognition with augmented synthesized speech,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 996–1002.
- [9] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: an ASR corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [10] Jin Xu, Xu Tan, Yi Ren, Tao Qin, Jian Li, Sheng Zhao, and Tie-Yan Liu, “LRSpeech: Extremely low-resource speech synthesis and recognition,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 2802–2812.
- [11] Roger Hsiao, Ralf Meermeier, Tim Ng, Zhongqiang Huang, Maxwell Jordan, Enoch Kan, Tanel Alumäe, Jan Silovsky, William Hartmann, Francis Keith, et al., “Sage: The new BBN speech processing platform,” in *Proc. INTERSPEECH 2016*, San Francisco, California, USA, 2016, pp. 3022–3026.
- [12] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, et al., “The Kaldi speech recognition toolkit,” in *IEEE 2011 workshop on automatic speech recognition and understanding*. IEEE Signal Processing Society, 2011, number CONF.
- [13] J.Z. Ma, F. Keith, T. Ng, M.-h. Siu, and O. Kimball, “Improving deliverable speech-to-text systems with multilingual knowledge transfer,” in *Proc. INTERSPEECH 2017*, Stockholm, Sweden, 2017, pp. 127–131.
- [14] Daniel Povey, Gaofeng Cheng, Yiming Wang, Ke Li, Hainan Xu, Mahsa Yarmohammadi, and Sanjeev Khudanpur, “Semi-orthogonal low-rank matrix factorization for deep neural networks.,” in *Interspeech*, 2018, pp. 3743–3747.
- [15] Alan W Black, “CLUSTERGEN: A statistical parametric synthesizer using trajectory modeling,” in *Ninth International Conference on Spoken Language Processing*, 2006.
- [16] Alan Black, Paul Taylor, and Richard Caley, “The festival speech synthesis system.”
- [17] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet, “Front-end factor analysis for speaker verification,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.
- [18] Ryan Prenger, Rafael Valle, and Bryan Catanzaro, “Waveglow: A flow-based generative network for speech synthesis,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 3617–3621.
- [19] Keith Ito and Linda Johnson, “The LJ speech dataset,” <https://keithito.com/LJ-Speech-Dataset>, 2017.
- [20] John Kominek and Alan W Black, “The cmu arctic speech databases,” in *Fifth JSCA workshop on speech synthesis*, 2004.
- [21] Merlijn Blaauw, Jordi Bonada, and Ryunosuke Daido, “Data efficient voice cloning for neural singing synthesis,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6840–6844.
- [22] John Kominek, Tanja Schultz, and Alan W Black, “Synthesizer voice quality of new languages calibrated with mean mel cepstral distortion,” in *Spoken Languages Technologies for Under-Resourced Languages*, 2008.