End-to-End Bangla Speech Synthesis

Prithwiraj Bhattacharjee, Rajan Saha Raju, Arif Ahmad, M. Shahidur Rahman Department of Computer Science and Engineering, Shahjalal University of Science & Technology, Sylhet, Bangladesh Email: {prithwiraj12, rajan10}@student.sust.edu, {arif_ahmad-cse, rahmanms}@sust.edu

Abstract—Text-to-Speech (TTS) system is a system where speech is synthesized from a given text following any particular approach. Concatenative synthesis, Hidden Markov Model (HMM) based synthesis, Deep Learning (DL) based synthesis with multiple building blocks, etc. are the main approaches for implementing a TTS system. Here, we are presenting our deep learning-based end-to-end Bangla speech synthesis system. It has been implemented with minimal human annotation using only 3 major components (Encoder, Decoder, Post-processing net including waveform synthesis). It does not require any frontend preprocessor and Grapheme-to-Phoneme (G2P) converter. Our model has been trained with phonetically balanced 20 hours of single speaker speech data. It has obtained a 3.79 Mean Opinion Score (MOS) on a scale of 5.0 as subjective evaluation and a 0.77 Perceptual Evaluation of Speech Quality(PESQ) score on a scale of [-0.5, 4.5] as objective evaluation. It is outperforming all existing non-commercial state-of-the-art Bangla TTS systems based on naturalness.

Index Terms—Text-to-Speech; End-to-End; MOS; PESQ; Tacotron; Griffin-lim;

I. INTRODUCTION

Text-to-Speech (TTS) system is a system that generates speech from any given text. Both Digital Signal Processing (DSP) and Natural Language Processing (NLP) are being used to develop such a system. In previous years, some great works have been done on TTS research. But many things should be taken under consideration for developing a better TTS system. If we preciously indicate language-specific TTS systems then there lies a vast area to work on. Bangla language is the fifthmost spoken language in the world. Technology like Bangla TTS system has a great impact on our daily life. But no significant implementation of Bangla TTS is available because all the existing systems are based on traditional approaches. So, we planned to take an initiative to develop a better Bangla TTS system by using the current state-of-the-art technologies.

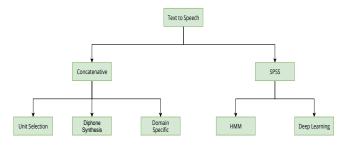


Figure 1: Classification of TTS system

TTS systems are mainly classified into two major parts. Concatenative TTS and Statistical Parametric Speech Syn-

thesis (SPSS) TTS. Some concatenative TTS systems are unit selection, diphone synthesis, and domain-specific synthesis. Unit selection requires pre-recorded speech and diphone synthesis needs all the sound-to-sound transitions or diphones. Hidden Markov Model (HMM) based TTS and Deep Neural Network (DNN) based TTS are SPSS TTS. Major components of traditional DNN based TTS systems are frontend preprocessing, acoustic modeling, duration modeling, and vocoder. In recent times, end-to-end TTS systems have become very popular. This neural network-based TTS system neither needs any pre-recorded speech nor any building blocks such as frontend preprocessing. Though researches on Bangla TTS started a long time back, we continued the legacy by doing researches on the latest TTS technologies. TTS system is very effective for physically impaired or blind people. Again modern business, airports, and all other notable sectors use the TTS system. These are some of our research motivations. Our discriminative research contributions are as follows:

- This is the first-ever Bangla end-to-end TTS system. It does not require complex linguistic and acoustic features as input.
- We built phonetically rich studio-quality speech data containing more than 40 hours of speech which more than any other publicly available dataset with the help of two professional voice artists.
- In terms of naturalness, the system outperforms among all other non-commercial Bangla TTS systems.

II. RELATED WORK

The very first attempt of implementing Bangla TTS system was concatenative developed by Alam et. al. [1]. It was a unit selection-based TTS system named Katha where they used festival¹ toolkit for front-end preprocessing. In 2009, a diphone concatenation-based TTS system Subachan was developed by Naser et. al. [2]. In 2014, Mukharjee et. al. [3] developed the first HMM-based SPSS system for Bangla. But over smoothness of acoustic modeling and limitations of vocoders resulted in producing less natural output in HMMbased system. In 1993, Weijters et. al. [4] first gave an idea of implementing neural networks for acoustic modeling. But when Hinton et. al. [5], ze et. al. [6] and Ling et. al. [7] individually implemented Deep neural networks (DNNs) for acoustic modeling, a great enhancement occurred in TTS research. Recurrent Neural Networks (RNNs) of Ze et. al. [6], Long Short-Term Memory RNNs (LSTM-RNNs) of Qian

http://www.cstr.ed.ac.uk/projects/festival/

et. al. [8], and deep Bidirectional LSTM-RNNs (BLSTM-RNNs) of Fan et. al. [9] significantly improved the performance of SPSS systems. Gutkin et. al. [10] from google released a commercial Bangla TTS system but the project is not publicly available. Although Google has released some language-specific resources on github² for enhancing Bangla TTS research. Potard et. al. [11] developed an open-source TTS system named Idlak Tangle taking help from another open-source speech recognition system named Kaldi proposed by Povey et. al. [12]. Wu et. al. [13] from Edinburgh University developed another DNN based open source TTS system named Merlin where they used ossian³ and festival frontend preprocessors and WORLD vocoder of Morise et. al. [14] for synthesizing speech. Taking help from these open source tools. Deka et. al. [15] developed a TTS system for a low resourced language. In 2019, Raju et. al. [16] also developed a DNN based Bangla SPSS TTS using this Merlin toolkit. In most cases, SPSS methods are better than concatenative approaches but they need a huge amount of human annotation such as language-specific text processing. So, researchers are now interested to develop end-to-end TTS systems where speech is directly generated from (text, speech) pairs without having any preprocessing. Inspiring from this concept, Arik et. al. [17], Gibiansky et. al. [18] and Ping et. al. [19] from Baidu developed Deep Voice 1, 2, and 3 respectively. Then Sotelo et. al. [20] from MILA developed Char2Wav and Ren et. al. [21] from Microsoft developed Fastspeech. Google has also published some TTS research-based papers named WaveNet by Oord et. al. [22], Tacotron by Wang et. al. [23] and Tacotron 2 by Shen et. al. [24]. Though commercial companies do not release their project publicly, individual researchers release their implementation online. By observing all these cases, we adopted a good implementation released by Keithito⁴. It is based on tacotron developed by Wang et. al. [23]. Thus we implemented our end-to-end Bangla TTS system.

III. DATASET COLLECTION AND PREPROCESSING

To build a good TTS system we need to collect significant speech data with corresponding transcriptions. Google has released 3 hours of Bangla speech data. Alam et. al. [25] from Brac University also released 13 hours of speech data. These are some notable publicly available datasets. But modern research stated that for developing a good TTS system we need at least 20 hours of speech data. So, we started to prepare a standard and notable dataset from scratch in our institution. After the preparation, we did some preprocessing to make it compatible with training the model.

A. Text Corpus

From Honnet et.al. [26], Sonobe et. al. [27], and Gab-drakhmanov et. al. [28] we got the concept of developing a phonetically balanced text corpus. We then collected our text data from various domains confirming that our corpus

has all possible punctuations for Bangla. Our final dataset contains more than 12500 utterances. A preview of our dataset is presented in table I.

Total sentences	12,537
Total words	1,22,627
Total unique words	24,582
Minimum words in a sentence	3
Maximum words in a sentence	20
Average words in a sentence	9.78
Total duration of speech (hours)	20:14:21
Average duration of each sentence (seconds)	5.81

TABLE I: Summary of Dataset

B. Speech Data

After completion of preparing the text corpus, we then started recording the speech data. In our university, a sound-proof lab has been built for recording the speech. Two professional voice artists (one male and one female) were employed to record speech for several consecutive hours. The sample rate of wave files is 48KHz. The duration of the collected speech is around 20 hours for each speaker. To make the data compatible with the tacotron model, we did some preprocessing. We removed all the sentences containing less than or equal to 3 words, greater than or equal to 12 words, and their corresponding speech. Then we removed the silence from starting and ending part from all the remaining waveforms. Thus we made our dataset compatible for training the model.

C. Text Normalization

Text normalization is the process of converting a raw text to its pronounceable form. To remove the ambiguities, nonstandard words should be converted into their standard pronunciation. Table II shows different examples of Bangla text normalization. So, we handled the issues like ambiguities of numerical words, abbreviations, etc., and normalized our full training data for generating more accurate pronunciation.

Raw Text	Category	Normalized Text
<i>৫৬</i> ১৫২৩	Phone Number	পাঁচ ছয় এক পাঁচ দুই তিন
	Amount	পাঁচ লক্ষ একষট্টি হাজার পাঁচশ তেইশ
5 2.60	Time	বারোটা পঞ্চাশ
	Fraction Money	বারো টাকা পঞ্চাশ পয়সা
১৬-১২-১৯৭১	Date	ষোলই ডিসেম্বর ঊনিশো একাত্তর
ড.	Abbreviation	ডক্টর

TABLE II: Bangla Text Normalization

IV. SYSTEM ARCHITECTURE

Wang et. al. [23] proposed that tacotron is an end-toend generative TTS model based on sequence-to-sequence model by Sutskever et. al. [29]. It also uses the attention mechanism of Bahdanau et. al. [30]. Tacotron uses Griffin-Lim reconstruction algorithm to generate waveform from spectrogram where Tacotron 2 [24] uses WaveNet vocoder to

²https://github.com/google/language-resources/

³https://github.com/CSTR-Edinburgh/Ossian

⁴https://github.com/keithito/tacotron

invert mel spectrogram into waveform. But WaveNet vocoder approach is slower than Griffin-Lim based synthesizer. So, we pick this model to develop our end-to-end Bangla TTS system. This system takes characters as input and generates a spectrogram as output. Tacotron TTS system does not need any frontend preprocessing and Grapheme-to-Phoneme (G2P) model. Training starts from scratch with random initialization. The system has three major components which drive the core methodology of the system. The components are given below:

- An Encoder
- An Attention-based Decoder
- A Post-processing Net

There are some other components like the CBHG module, Griffin-lim reconstruction, and pre-net. We will now describe these components.

A. CBHG Module

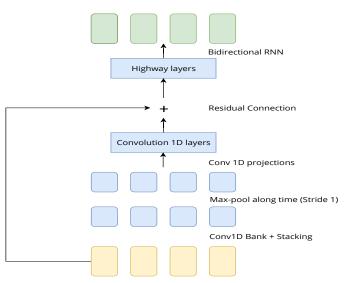


Figure 2: CBHG Module from [23, p. 3]

CBHG stands for a 1-D Convolutional Bank, Highway Networks, and a Bidirectional Gated Recurrent Unit (GRU). In CBHG, the input sequence has been passed through 1-D convolution filters. After doing some processing the convoluted output is then added with the original input sequence by residual connections introduced by He et. al. [31]. Then the output is feed into a multi-layer highway network which is used for extracting high-level features. Lastly, a Bidirectional GRU is used at the top level for extracting sequential features. By using both forward and backward context the GRU completes the extraction. Both batch normalization of Ioffe et. al. [32] and residual connections etc. improved the generalization. Figure 2 shows the CBHG module.

B. Major Components

Tacotron has three major building blocks. An encoder, an attention-based decoder, and a post-processing net. The complete system architecture is shown in figure 3.

1) Encoder: Sequential representations of text are extracted by an encoder. Character sequences are the inputs for the encoder represented as a one-hot vector. These one-hot vectors are then embedded into a continuous vector representation. For all this to happen, some other components are also being used in the encoder part.

- Pre-net
- · CBHG Module

Pre-net is a combination of some non-linear transformations and applied to each embedding. It helps to improve the generalization as well as the convergence of the system. CBHG module transforms the output to a final representation and an attention module uses the modified output for further processing. CBHG-based encoder reduces the overfitting and mispronunciation of words.

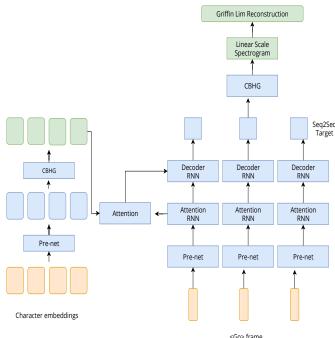


Figure 3: System Architecture from [23, p. 2]

- 2) Attention Based Decoder: Vinyals et. al. [33] attention based decoder is used in our system. The decoder part includes the following things:
 - Attention RNN
 - Pre-net
 - Decoder RNN
 - Sequence-to-Sequence Target

At first, the encoder output is concatenated with the attention RNN and fed into the decoder RNN. An initial pre-net is also fed into the decoder with all zero frames. Decoder consists of a stack of GRUs connected via residual connections. The sequence-to-sequence target model receives the output of the decoder RNN. Then a fixed group of sequence-to-sequence targets started feeding into a CBHG module in the post-processing net part.

- 3) Post Processing Net: Speech is synthesized in the post-processing net. The post-processing net knows how to predict spectral magnitude based on a linear-frequency scale. It has both forward and backward information and can observe the decoded sequence of target output. Two main sub-components play a vital role in the post-processing net.
 - · CBHG Module
 - Griffin-Lim algorithm

CBHG module takes the sequence-to-sequence target output as input and produces the linear-scale spectrogram. Griffin-Lim algorithm, introduced by Griffin et. al. [34] is used for generating speech from the spectrogram. This is called spectrogram inversion.

4) Griffin-Lim Reconstruction: This algorithm is by far the fastest spectrogram inversion algorithm. It generates speech from the raw spectrogram. Abadi et. al. [35] implemented the Griffin-Lim algorithm with TensorFlow and we followed this procedure. Though this algorithm converges after 50 iterations, we went up to 60 iterations for more accuracy.

V. OUR SYSTEM

The system has been trained for several days and generates a new model checkpoint after every 1000 iterations. Each checkpoint synthesizes a waveform and generates a graph of attention alignments between encoder and decoder timestamp with loss information. Figure 4 shows the attention alignment graph of our model. The best checkpoint has been found after 228k steps of training with a loss of 6.017%.

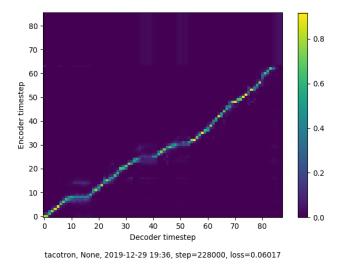


Figure 4: Attention alignments on a test phrase after 228k steps

The Y-axis of this graph denotes the encoder timestamp which at each step took an input character as well as its current state and gave an output of 100 real vectors representing the status of each moment. The X-axis denotes the decoder timestamp which reads all the status vectors of the encoder and step-by-step generated audio frames or mel-spectrogram. The

diagonal line represents the result of those audio frames based on each input character. The more the graph is diagonally aligned the more the model performs well.

VI. RESULTS

There are mainly two types of evaluation that exist to measure the performance of a TTS system. **Objective Evaluation** and **Subjective Evaluation**. From these evaluations, we will now show the numerical results as well as the graphical comparisons of different existing Bangla TTS systems.

A. Objective Evaluation

Objective evaluation is nothing but a mathematical comparison of the generated signal and the original signal. Perceptual Evaluation of Speech Quality (PESQ) [36] is one of the most adopted metrics for objective evaluation.

1) PESQ: There are different types of PESQ measurement exist. MOS-LQO and raw-PESQ are mostly being used. The ranges of raw-PESQ and MOS-LQO are [-0.5,4.5] and [1.0, 5.0] respectively. We measured the raw-PESQ score for evaluating our model performance. As both the original and generated waveforms are required to calculate PESQ, we took 100 random sentences and their original waveforms from the test dataset. From those 100 sentences, we synthesized the waveforms using our TTS system. After that, the original and generated waveforms of corresponding sentences were sent to the PESQ system simultaneously. Figure 5 shows the PESQ score of those 100 generated waveforms. We then calculated the average PESQ score of those 100 waveforms and got 0.77 as the average PESQ score.

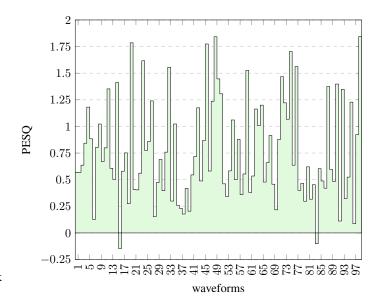


Figure 5: PESQ score of 100 synthesized wave

B. Subjective Evaluation

Subjective evaluation is the user-oriented evaluation. Mean Opinion Score (MOS) [37] is one of the most popular metrics for subjective evaluation and we also calculated the MOS to evaluate our TTS system.

1) MOS: MOS measured the naturalness of a TTS system. Five labels are given to some users for scoring the system. These five labels are Excellent, Good, Fair, Poor, and Bad. After listening to the waveforms, the individual user chooses a label to rate each of the synthesized waveforms. These five

Label	Rating
Excellent	5
Good	4
Fair	3
Poor	2
Bad	1

TABLE III: Label vs Rating

labels are mapped with some pre-defined numbers. Table III shows the mapping between these five labels and the numbers. After that, the arithmetic mean of those numbers has been calculated using the equation 1 which we called MOS.

$$MOS = \frac{1}{N} \sum_{n=1}^{N} R_n \tag{1}$$

Individual rating is an integer number but the final result can be a real value. For rating our system, 65 native Bengali speakers from different backgrounds participated. We provided 10 waveforms to them generated by our TTS system. Those 10 waveforms are available at this Github⁵ link. They gave a label-based rating to those 10 waveforms based on naturalness. So, we received 65 ratings for each of the 10 waveforms and calculated an average rating of the individual 10 waveforms. Finally, we calculated the arithmetic mean of those 10 individual ratings using the formula 1 and got 3.79 as Mean Opinion Score (MOS). Figure 6 shows the individual MOS of 10 waveforms.

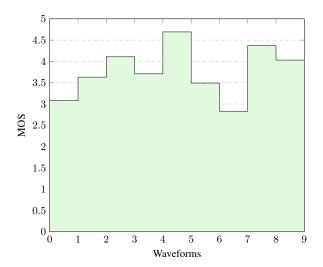


Figure 6: MOS of 10 synthesized wave

C. Discussion and Analysis

We are now going to show the graphical comparisons of various existing Bangla TTS systems. We calculated PESQ

and MOS for our Tacotron TTS as well as for Google Bangla TTS [10], SPSS TTS [16], and Subachan TTS [2]. The bar chart of figure 7 shows the comparison of PESQ scores of different TTS systems. Subachan TTS and SPSS TTS have obtained 0.45 and 0.53 PESQ scores respectively. Our current Tacotron TTS has achieved a 0.77 PESQ score which is better than noncommercial SPSS TTS and Subachan TTS. Here we can observe a gradual advancement of our system. Another comparison of various TTS systems is shown in figure 8 concerning Mean Opinion Score (MOS). We have achieved 3.79 MOS which is better than the previously implemented systems like Subachan TTS and SPSS TTS. There are no significant works on Bangla TTS using deep learning. A very recent work Byakto Speech [38] has obtained a 3.23 MOS score. In some cases, it is even comparable with commercially implemented Google Bangla TTS.

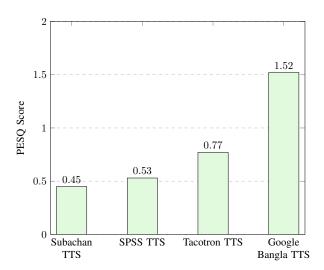


Figure 7: PESQ Scores of Different TTS Systems

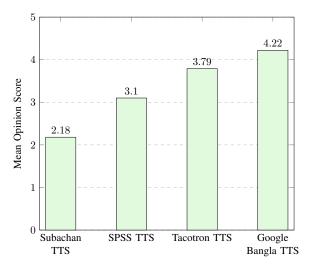


Figure 8: MOS of Different TTS Systems

⁵https://sustbn.github.io/

VII. CONCLUSION

Tacotron TTS or end-to-end Bangla TTS is the current state-of-the-art TTS system. It overcomes the drawbacks of the traditional Deep Neural Network (DNN) based TTS system which needs more building blocks, hand-engineered features, etc. Based on naturalness, achieving a 3.79 score as MOS is very significant. This score outperforms all existing non-commercial Bangla TTS systems. Moreover, based on the objective evaluation we got a 0.77 PESQ score which also correlates with the MOS. We used the Griffin-lim spectrogram inversion algorithm which is till now the fastest one. Though text normalization has been done before feeding the data into the network, modern research is proposing the automation of text normalization.

REFERENCES

- F. Alam, P. K. Nath, and M. Khan, "Text-to-speech for bangla language using festival," 2007.
- [2] A. Naser, D. Aich, and M. R. Amin, "Implementation of subachan: Bengali text-to-speech synthesis software," in *International Conference on Electrical & Computer Engineering (ICECE 2010)*. IEEE, 2010, pp. 574–577.
- [3] S. Mukherjee and S. K. D. Mandal, "A bengali hmm based speech synthesis system," arXiv preprint arXiv:1406.3915, 2014.
- [4] T. Weijters and J. Thole, "Speech synthesis with artificial neural networks," in *IEEE International Conference on Neural Networks*. IEEE, 1993, pp. 1764–1769.
- [5] G. Hinton, L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal processing magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [6] H. Ze, A. Senior, and M. Schuster, "Statistical parametric speech synthesis using deep neural networks," in 2013 ieee international conference on acoustics, speech and signal processing. IEEE, 2013, pp. 7962–7966
- [7] Z.-H. Ling, S.-Y. Kang, H. Zen, A. Senior, M. Schuster, X.-J. Qian, H. M. Meng, and L. Deng, "Deep learning for acoustic modeling in parametric speech generation: A systematic review of existing techniques and tuture trends," *IEEE Signal Processing Magazine*, vol. 32, no. 3, pp. 35–52, 2015.
- [8] Y. Qian, Y. Fan, W. Hu, and F. K. Soong, "On the training aspects of deep neural network (dnn) for parametric tts synthesis," in 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2014, pp. 3829–3833.
- [9] Y. Fan, Y. Qian, F.-L. Xie, and F. K. Soong, "Tts synthesis with bidirectional lstm based recurrent neural networks," in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [10] A. Gutkin, L. Ha, M. Jansche, K. Pipatsrisawat, and R. Sproat, "Tts for low resource languages: A bangla synthesizer," in *Proceedings of the Tenth International Conference on Language Resources and Evaluation* (LREC 2016), 2016, pp. 2005–2010.
- [11] B. Potard, M. P. Aylett, D. A. Baude, and P. Motlicek, "Idlak tangle: An open source kaldi based parametric speech synthesiser based on dnn." in *INTERSPEECH*, 2016, pp. 2293–2297.
- [12] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz et al., "The kaldi speech recognition toolkit," in *IEEE 2011 workshop on automatic speech recognition and understanding*, no. CONF. IEEE Signal Processing Society, 2011.
- [13] Z. Wu, O. Watts, and S. King, "Merlin: An open source neural network speech synthesis system." in SSW, 2016, pp. 202–207.
- [14] M. Morise, F. Yokomori, and K. Ozawa, "World: a vocoder-based high-quality speech synthesis system for real-time applications," *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877–1884, 2016.

- [15] A. Deka, P. Sarmah, K. Samudravijaya, and S. Prasanna, "Development of assamese text-to-speech system using deep neural network," in 2019 National Conference on Communications (NCC). IEEE, 2019, pp. 1–5.
- [16] R. S. Raju, P. Bhattacharjee, A. Ahmad, and M. S. Rahman, "A bangla text-to-speech system using deep neural networks," in 2019 International Conference on Bangla Speech and Language Processing (ICBSLP). IEEE, 2019, pp. 1–5.
- [17] S. Ö. Arik, M. Chrzanowski, A. Coates, G. Diamos, A. Gibiansky, Y. Kang, X. Li, J. Miller, A. Ng, J. Raiman et al., "Deep voice: Realtime neural text-to-speech," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 195–204.
- [18] A. Gibiansky, S. Arik, G. Diamos, J. Miller, K. Peng, W. Ping, J. Raiman, and Y. Zhou, "Deep voice 2: Multi-speaker neural text-tospeech," in *Advances in neural information processing systems*, 2017, pp. 2962–2970.
- [19] W. Ping, K. Peng, A. Gibiansky, S. O. Arik, A. Kannan, S. Narang, J. Raiman, and J. Miller, "Deep voice 3: Scaling text-to-speech with convolutional sequence learning," arXiv preprint arXiv:1710.07654, 2017.
- [20] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. Courville, and Y. Bengio, "Char2wav: End-to-end speech synthesis," 2017.
- [21] Y. Ren, Y. Ruan, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, "Fastspeech: Fast, robust and controllable text to speech," arXiv preprint arXiv:1905.09263, 2019.
- [22] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," arXiv preprint arXiv:1609.03499, 2016.
- [23] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio et al., "Tacotron: Towards end-to-end speech synthesis," arXiv preprint arXiv:1703.10135, 2017.
- [24] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. 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.
- [25] F. Alam, S. M. Habib, D. A. Sultana, and M. Khan, "Development of annotated bangla speech corpora," in *Spoken Languages Technologies* for Under-Resourced Languages, 2010.
- [26] P.-E. Honnet, A. Lazaridis, P. N. Garner, and J. Yamagishi, "The siwis french speech synthesis database? design and recording of a high quality french database for speech synthesis," Idiap, Tech. Rep., 2017.
- [27] R. Sonobe, S. Takamichi, and H. Saruwatari, "Jsut corpus: free large-scale japanese speech corpus for end-to-end speech synthesis," arXiv preprint arXiv:1711.00354, 2017.
- [28] L. Gabdrakhmanov, R. Garaev, and E. Razinkov, "Ruslan: Russian spoken language corpus for speech synthesis," arXiv preprint arXiv:1906.11645, 2019.
- [29] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing* systems, 2014, pp. 3104–3112.
- [30] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," arXiv preprint arXiv:1409.0473, 2014.
- [31] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision* and pattern recognition, 2016, pp. 770–778.
- [32] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," arXiv preprint arXiv:1502.03167, 2015.
- [33] O. Vinyals, Ł. Kaiser, T. Koo, S. Petrov, I. Sutskever, and G. Hinton, "Grammar as a foreign language," in *Advances in neural information processing systems*, 2015, pp. 2773–2781.
- [34] D. Griffin and J. Lim, "Signal estimation from modified short-time fourier transform," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 32, no. 2, pp. 236–243, 1984.
- [35] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.
- [36] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in 2001 IEEE International

- Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No. 01CH37221), vol. 2. IEEE, 2001, pp. 749–752.
- [37] M. Viswanathan and M. Viswanathan, "Measuring speech quality for text-to-speech systems: development and assessment of a modified mean opinion score (mos) scale," *Computer Speech & Language*, vol. 19, no. 1, pp. 55–83, 2005.
- no. 1, pp. 55–83, 2005.

 [38] Z. A. Nazi and S. M. T. Huda, "Byakto speech: Real-time long speech synthesis with convolutional neural network: Transfer learning from english to bangla," arXiv preprint arXiv:2106.03937, 2021.