

Improved Child Text-to-Speech Synthesis through Fastpitch-based Transfer Learning

Rishabh Jain
C3 Imaging Research Center
University of Galway
Galway, Ireland
rishabh.jain@universityofgalway.ie

Peter Corcoran
C3 Imaging Research Center
University of Galway
Galway, Ireland
peter.corcoran@universityofgalway.ie

Abstract— Speech synthesis technology has witnessed significant advancements in recent years, enabling the creation of natural and expressive synthetic speech. One area of particular interest is the generation of synthetic child speech, which presents unique challenges due to children's distinct vocal characteristics and developmental stages. This paper presents a novel approach that leverages the Fastpitch text-to-speech (TTS) model for generating high-quality synthetic child speech. This study uses the transfer learning training pipeline. The approach involved finetuning a multi-speaker TTS model to work with child speech. We use the 'cleaned' version of the publicly available MyST dataset (55 hours) for our finetuning experiments. We also release a prototype dataset of synthetic speech samples generated from this research together with model code to support further research. By using a pretrained MOSNet, we conducted an objective assessment that showed a significant correlation between real and synthetic child voices. Additionally, to validate the intelligibility of the generated speech, we employed an automatic speech recognition (ASR) model to compare the word error rates (WER) of real and synthetic child voices. The speaker similarity between the real and generated speech is also measured using a pretrained speaker encoder.

Keywords—Fastpitch, synthetic speech, child speech, wav2vec2, MOSNet, Wavglow, MyST dataset.

I. INTRODUCTION

Speech synthesis technology has witnessed significant advancements in recent years, enabling the creation of natural and expressive synthetic speech. One area of particular interest is the generation of synthetic child speech, which presents unique challenges due to children's distinct vocal characteristics and developmental stages. Early research on Text-to-Speech (TTS) synthesis began several decades ago, primarily using concatenative and parametric methods [1]–[4]. While these methods generated speech from text, the resulting audio lacked naturalness and sounded robotic. Recent advancements in TTS models, mainly based on deep neural networks (DNN), have significantly improved the quality of synthesized speech. Tacotron [5], a neural sequence-to-sequence model, marked a notable improvement in speech synthesis quality. Subsequent models like Tacotron2 [6], FastSpeech [7], FastSpeech2 [8], FlowTTS [9], GlowTTS [10], Fastpitch [11], and Adaspeech [12] have further evolved TTS capabilities and improved TTS speech quality. Deepvoice2 [13] introduced the use of speaker verification models [14]–[16] to achieve multi-speaker TTS [17]–[21].

Child-TTS (CTTS), or TTS synthesis for child speech is currently limited due to the scarcity of child voice datasets and the challenges associated with their creation. Collecting child speech data for TTS research is challenging. Most TTS datasets are created in studios with expensive equipment,

tailored for adult voices. While the pitch for adults typically falls between 70 to 250 Hz, children's speech ranges from 200 to 500 Hz [22]. Additionally, child speech exhibits distinct characteristics from adult speech, such as a higher fundamental frequency and variable speaking rates compared to adults [23]–[26]. Moreover, children tend to have longer phoneme durations and different prosody features due to their smaller vocal tracts [27]–[29].

This research aims to harness the potential of state-of-the-art (SOTA) TTS methods such as Fastpitch [11] to construct a pipeline for synthesizing children's voices while minimizing data requirements. The primary objective is to demonstrate the pipeline's ability to reliably generate a variety of self-consistent, distinct children's voices. Fastpitch utilizes a pitch prediction and duration prediction module which captures pitch variations in speech and enables more precise control over the speaking rate. This study uses an existing multispeaker children's speech dataset [30], which was cleaned to make it more suitable for CTTS research [31]. Subsequently, Fastpitch was trained on the cleaned dataset to generate synthetic speech for multiple child speakers, serving as a proof of concept.

By incorporating Fastpitch into the synthesis pipeline, we can effectively capture the unique prosodic features and intonation patterns present in speech. Our objective is to further optimize this model for child speech to accommodate individual characteristics, such as gender and regional accents, to produce realistic synthetic child voices. By using this approach, we intend to overcome the limitations of traditional TTS systems that often fail to capture the naturalness and authenticity of child-like speech. Our hypothesis is centered on the idea that pretraining the TTS model on adult speech data and subsequently finetuning it with child speech data can facilitate the synthesis of artificial child speech.

As part of this research, we also release a small set of synthetic datasets generated from this research. Objective evaluations were conducted on the synthesized child voices, comparing them to real child voices in terms of various acoustic features and Mean Opinion Score (MOS). The evaluation encompassed factors such as 'Naturalness', 'Intelligibility', and 'Speaker Similarity'. Furthermore, we compared this approach with our previously reported Tacotron 2 TTS pipeline for the child speech synthesis [32]. In this study, no subjective evaluation was conducted; however, it will be taken into consideration for future research.

The potential applications of this research are wide-ranging and impactful such as educational tools, audiobooks

for children, language learning, interactive games and toys, virtual learning companions, and child-friendly voice assistants and chatbots to name a few. Such a pipeline would also enable the creation of large synthetic datasets, which could, in turn, enhance other areas of child speech research, such as speaker recognition and automatic speech recognition [33], [34].

II. METHODOLOGY

Fastpitch is a fully parallel TTS model conditioned on fundamental frequency contours. By incorporating Fastpitch into the synthesis pipeline, we can effectively capture the unique prosodic features and intonation patterns present in child speech. We present a multispeaker framework for TTS using a transfer learning approach that uses Waveglow vocoder for audio synthesis. We also evaluate this methodology using different objective evaluation methods to provide the validity of this approach. Fastpitch is used in this work due to its various advantages such as faster inference speed, improved prosody control, enhanced naturalness, duration control, multilingual support, and simplified architecture as compared to previous TTS approaches.

A. Datasets

In this section, we give an overview of the datasets used to finetune our pipeline and to implement some of our evaluation methods.

1) TTS Datasets:

These datasets are used for the TTS experiments for pretraining and finetuning the TTS model.

LibriTTS [35]: The LibriTTS corpus is an adult speech dataset that includes 585 hours of speech data sampled at a rate of 24kHz, obtained from a diverse set of 2,456 speakers. LibriTTS is widely used in research for training and evaluating text-to-speech systems.

MyST [30]: My Science Tutor (MyST) Corpus [36] is an American English child speech dataset from 1371 students containing over 393 hours of audio data out of which 197 hours are fully transcribed. We use the cleaned version of this dataset (derived from [31]), with 65 hours of speech divided into two subsets: 55 hours for training, called ‘MyST_train’ and 10 hours for testing, called ‘MyST_test’. This 55 hours of training data is used for TTS training.

2) Text Datasets:

These datasets are used during inference as input text for the TTS model to generate data samples from the finetuned synthetic child voices.

Harvard Sentences [36]: Harvard sentences consist of 720 sentences that are carefully designed to be phonetically balanced. These sentences effectively encompass a wide range of phonemes.

LJ Speech Sentences [37]: This dataset contains 13,100 sentences extracted from the LJ Speech dataset.

B. Multispeaker child TTS using Fastpitch

1) Fastpitch (Acoustic Model) [11]:

FastPitch is a streamlined TTS model with a simplified encoder-decoder architecture, designed for faster inference and improved prosody control. In the multispeaker FastPitch

TTS model, the input text is encoded using an encoder module, which typically comprises stacked layers of convolutional neural networks (CNNs) or recurrent neural networks (RNNs). The encoder processes the linguistic features of the text, such as phonemes or graphemes, and generates intermediate representations. The duration predictor module takes the intermediate representations from the encoder and predicts the duration of each phoneme or character in the input text. This enables the model to capture and generate natural speech rhythm and timing. The pitch predictor module takes the intermediate representations and predicts the fundamental frequency (F0) contour, controlling the pitch variations in the synthesized speech. The architecture of Fastpitch is detailed in Figure 1.

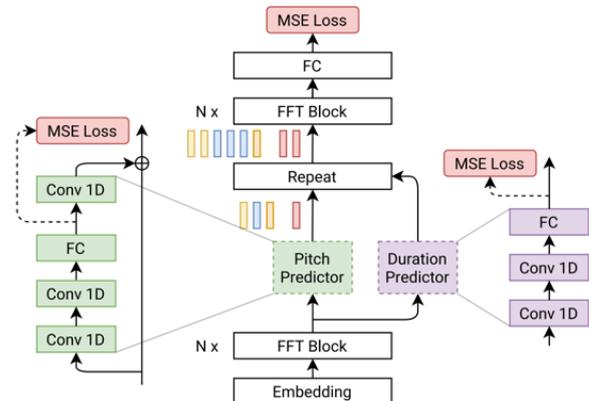


Fig. 1. Fastpitch Architecture [11].

We also condition the model on the speaker by adding a global speaker embedding [38] to the input tokens. The speaker embedding integration with the TTS framework [39]–[41] allows the model to capture the unique characteristics of different speakers. These embeddings encode speaker-related information in a vectorial representation for each speaker. During training, the model learns to associate speaker embeddings with the corresponding speakers, allowing it to generate speech that not only follows the desired linguistic content but also reflects the distinct vocal attributes of specific speakers. The primary loss function is the mean squared error (MSE) between the predicted mel-spectrogram and the target mel-spectrogram. Our work uses a newer version of Fastpitch, which is based on using the self-attention framework proposed in [38]. This enables the TTS model to learn speech-to-text alignment in parallel to TTS training instead of relying on an external aligner.

2) Transfer Learning Pipeline:

The proposed methodology involves pretraining the Fastpitch TTS model on a diverse dataset of adult speech, covering various age groups, linguistic backgrounds, and speech contexts. The LibriTTS dataset was used in this work. By finetuning the pretrained model on a smaller subset of the child speech dataset, such as MyST, will enable the model to learn the distinctive acoustic properties and pitch contours specific to child speech. Moreover, the model can be further optimized to accommodate individual characteristics, such as gender and regional accents, to synthesize more realistic CTTS voices.

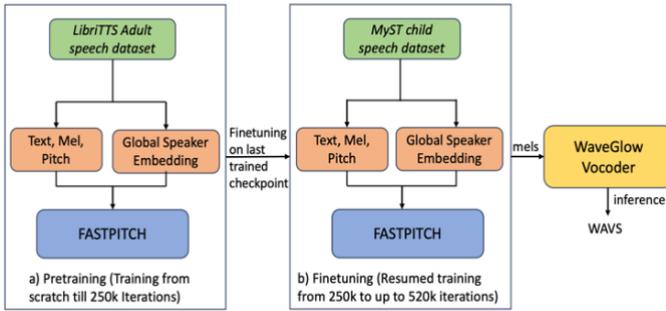


Fig. 2. Transfer learning pipeline: a) Pretraining: model being trained with LibriTTS dataset for up to 250k iterations. b) Finetuning: Resuming the acoustic model training with the MyST dataset from 250k iteration onwards up to 520k iteration.

The finetuning pipeline is kept consistent with our previous approach using Tacotron 2 [32] to allow for comparisons. Figure 2 describes the transfer learning pipeline. The model is first trained with the LibriTTS dataset (585 hours) for up to 250k iterations until a consistent low loss threshold is achieved, and the model starts to converge. After that, the model was finetuned for up to 520k additional steps using the MyST dataset (55 hours).

3) Waveglow (Vocoder) [42]:

WaveGlow is a SOTA vocoder model that generates high-quality and natural-sounding speech waveforms. It is based on a generative flow-based model architecture which models the distribution of speech waveforms. WaveGlow operates by taking a spectrogram representation of the speech as input and generating the corresponding waveform. The model employs an invertible neural network to transform the spectrogram into a latent space representation and then uses a series of invertible coupling layers to map this latent representation back into the waveform domain. Our WaveGlow model is trained on LibriTTS adult speech data which learns the complex relationship between spectrograms and waveforms. It was observed that Glow models [43]–[46] has popularly been used as a universal vocoder [45] and has been shown to work well with unseen speakers in multi-speaker models as well [47], [48]. Therefore, for the scope of this paper, WaveGlow (trained on LibriTTS) is used as a universal vocoder with synthetic child voices.

III. EXPERIMENTS

A. Training details

The implementation is obtained from Nvidia’s FastPitch Github¹. For our training and finetuning process, we utilized two A6000 40GB GPUs. We employed a learning rate of 0.1 and a weight decay factor of $1e-6$, maintaining consistency with their original implementation [11]. Additionally, the remaining hyperparameters were retained as per the provided implementation details. To ensure a smooth training process, we incorporated a warmup training step with a factor of 2000.

B. Experiments

1) Initial Experiments: These experiments involved using the LJ speech dataset for single-speaker finetuning. The model was first trained with LJ Speech and then finetuned

with a single speaker from the MyST dataset. The output audio obtained was quite noisy. We also tried training the LJ speech single-speaker dataset and finetuning it with the complete MyST dataset (considering it as a single-speaker dataset). However, the results obtained didn’t sound like child speech. Hence, finetuning on a single speaker was not explored further.

2) Main Experiments: These experiments involve multispeaker TTS training. The model was first trained with the LibriTTS dataset. Figure 3 shows an example loss curve of the LibriTTS training.

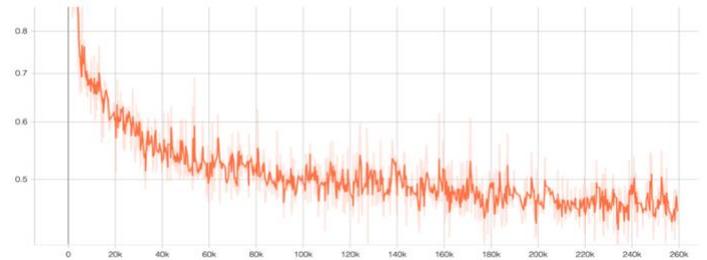


Fig. 3. LibriTTS pretraining curve (MSE loss vs. number of epochs)

It can be observed that for the first 2000 warmup steps, loss decreases gradually. After that loss decreases steadily until it reaches an average loss of 0.3 around 250k epoch. Since there was no improvement in loss function after that, it was decided to pause the training for further finetuning.

Further finetuning was performed from epoch 250k onwards on the MyST dataset. The loss increases until it starts to decrease around 260k epoch. From this point, there is a gradual decrease in loss until 520k steps. No significant improvement was observed in loss after this epoch. This was also verified by manually listening to generated audio files at an interval of every 50k epoch. After 550k epochs, the model exhibited signs of overfitting and began learning noise features in the MyST dataset, resulting in a decline in the quality of the synthesized audio.

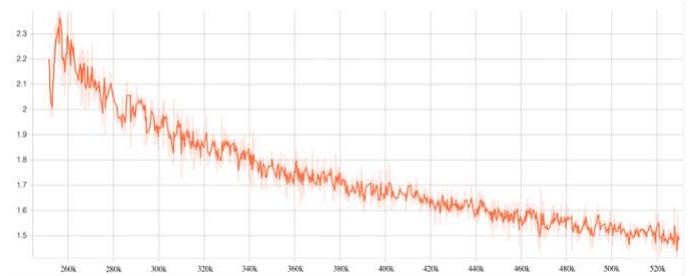


Fig. 4. MyST Finetuning curve (MSE loss vs. number of epochs).

C. Synthetic Datasets

We have generated two sets of synthetic child speech datasets. The dataset demographic is detailed in Table III. The dataset is made available through our Github². Since the dataset was generated at a 22Khz sampling rate, FFmpeg was used to convert the data into a 16khz sampling rate for objective evaluation. The dataset is made available in both sampling rates. The dataset details are available below:

¹ <https://github.com/NVIDIA/DeepLearningExamples/tree/master/PyTorch/SpeechSynthesis/FastPitch>

² https://github.com/C3Imaging/child_tts_fastpitch/

TABLE I. SYNTHETIC DATASET DEMOGRAPHICS

Dataset	Speakers	Hours	Utterances	data/speaker
CS_HS	40	29.02	28,800	43.53 minutes
CS_LJ	2	47.61	26,200	23.8 hours

1) *CS_HS* – This dataset used Harvard Sentences as a text reference to generate the synthetic child speech dataset. We selected the 40 speakers with the most amount of data in hours, from the LibriTTS dataset which was used to generate 40 child speakers. See Table 1 for more details.

2) *CS_LJ* – This dataset used LJ Speech transcripts as a text reference to generate the synthetic child speech dataset. We selected one male and one female speaker from the LibriTTS dataset which contained the most amount of training. These speakers were subjected to generate the child’s speech. See Table 1 for more details.

IV. RESULTS AND EVALUATION

Our experimental findings demonstrate the successful synthesis of child voices using our proposed methodology. To assess the validity of the generated speech, we conducted objective evaluations, specifically focusing on the aspects of Naturalness, Intelligibility, and Speaker similarity. Furthermore, we conducted a comparative analysis with our previous research, which involves synthesizing child speech using the Tacotron 2 model. For the evaluation process, we randomly selected 120 utterances from the original MyST dataset, Tacotron-based synthetic dataset, and Fastpitch-generated synthetic utterances (from III.B). This allowed us to systematically compare the quality of speech generated by both the Tacotron 2 [32] and Fastpitch models within the context of child speech synthesis.

A. Objective Naturalness Evaluation using the pretrained MOSNet [49]

TABLE II. MOSNET OUTPUT FOR 120 SAMPLES WITH 95% CONFIDENCE INTERVAL

Dataset	MOS
Adult speech (Librispeech test_clean)	3.78 ± 0.07
Original Child Speech [MyST]	2.91 ± 0.07
Tacotron 2 based synthetic child speech [32]	2.60 ± 0.06
Fastpitch based synthetic child speech [Our work]	3.10 ± 0.12

Table 1 provides the Mean Opinion Scores (MOS) for 120 different speech samples using the pretrained MOSNet model [49]. MOSNet, trained on adult speech, exhibits a high correlation with human MOS ratings. However, its generalization to child speech is doubtful. Therefore, we only use MOSNet in this study to explore the correlation between reference child audio and synthetic child audio. It acts as a measure to validate the ‘Naturalness’ of the speech. The original child speech from the MyST dataset received an average MOS of 2.91 ± 0.07 , indicating moderate acceptability. The Fastpitch-generated child speech indicates

higher quality than both the original speech and Tacotron2. These results suggest that the Fastpitch model, as implemented in our research, produces a strong correlation between synthetic child speech and real child speech.

B. Objective Intelligibility Evaluation using a pretrained wav2vec2 ASR System [50]

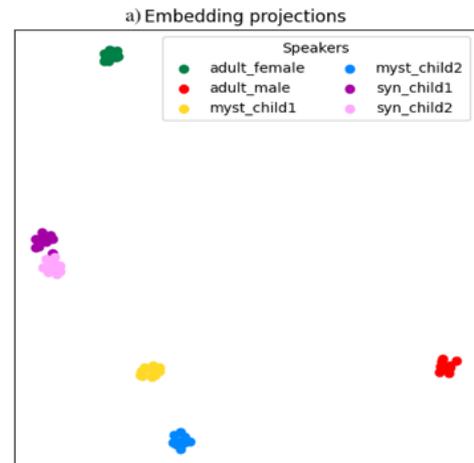
TABLE III. WER ON 120 RANDOMLY SELECTED UTTERANCES FROM ADULT SPEECH, REAL CHILD SPEECH, AND SYNTHETIC CHILD SPEECH USING THE WAV2VEC2 BASE ASR MODEL

Dataset	WER
Adult Speech (Librispeech test_clean)	3.43
Original Child Speech [MyST]	15.27
Tacotron 2 based synthetic child speech [32]	25.63
Fastpitch based synthetic child speech [Our work]	17.61

In this study, we employed the wav2vec2 base model³, which was finetuned with 960 hours of the Librispeech dataset, to evaluate the ‘Intelligibility’ of the generated child speech. Since wav2vec2 is a SOTA ASR model, it was intended to use this as a validity metric for the synthetic speech. Additionally, we conducted a comparative analysis with our previous approach utilizing Tacotron 2 [32]. Table II provides WER for different speech datasets. The adult speech dataset achieved a strong WER of 3.43, considering the model’s training on adult speech data. Our Fastpitch-based approach achieved a WER of 17.61, closely resembling the WER of the original child speech from the MyST dataset. Moreover, it surpassed the WER of the Tacotron 2 generated child speech, indicating improved performance over the synthetic child speech.

C. Speaker similarity verification using a pretrained speaker verification system [15]

Speaker similarity between a synthesized speech and a real speech can be calculated using a speaker verification system [15]. The pretrained speaker encoder from Resemblyzer⁴ was used to extract and visualize the speaker embeddings. This tool uses cosine distance to calculate the similarity between the two embeddings.



³ <https://github.com/facebookresearch/fairseq/tree/main/examples/wav2vec>

⁴ <https://github.com/resemble-ai/Resemblyzer>

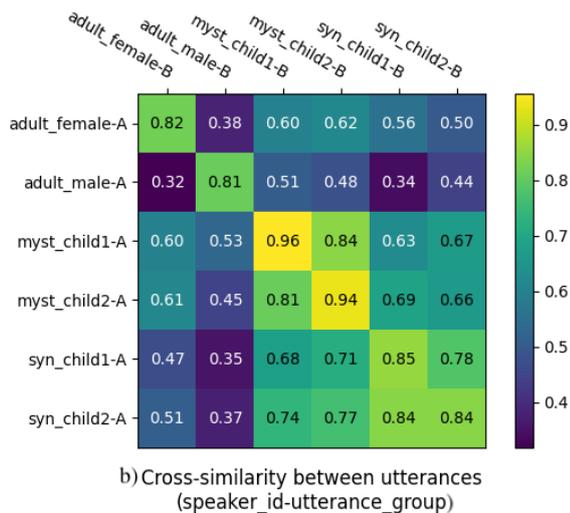


Fig. 5. a) Projections of embeddings between different real and synthetic child speech in comparison to adult speakers. b) Cross-similarity between 10 speakers in Set A and Set B.

For this evaluation, 6 speakers were randomly selected: 2 from LibriTTS [one male and one female], 2 from the MyST dataset, and 2 from the synthetically generated CS_HS dataset. We selected 10 utterances for each speaker in a random order. A visualization of this similarity in a 2D projection can be seen in Figure 5a. It can be observed that most of the child speakers (both real and synthetic) are very close in a cluster compared to adult male and female speakers.

To further demonstrate the similarity between real child speech and synthetic child speech, cosine similarity was used to calculate the cross-similarity between each speaker. All 6 speakers were divided into two sets, A and B. Speaker embeddings are extracted for each of the utterances for each of the sets and averaged together for each speaker. This gave us 6 unique speaker embeddings in sets A and B for each of the 6 speakers. Cosine similarity is finally used to measure the similarity between sets A and B. Figure 5b shows the plot for the cross similarity between 6 speakers. The similarity for most of the child and adult speech is between 0.34-0.53 whereas the similarity for synthetic child speech and real child speech is between 0.63-0.98. The average similarity between synthetic and real child voices is 77%. Hence, we can conclude that our synthetically generated child speech is quite close to real child speech in terms of speaker similarity.

V. CONCLUSION AND FUTURE WORK

This paper presents a pipeline for synthesizing child speech in scenarios with limited training data. The proposed approach involves cleaning an existing child speech dataset to create a small, curated dataset suitable for TTS training. A transfer learning technique is employed, utilizing pretraining on adult speech data and finetuning on child speech data. Objective evaluations using MOSNet demonstrate a strong correlation between real and synthesized child voices. Using a pretrained adult speech wav2vec2 ASR model, the WER for synthetic child voices was measured at 17.61, compared to a WER of 15.27 for real child voices. Speaker similarity evaluation using a pretrained speaker encoder yields an

average cosine similarity of 77% between synthetic speech and the original speakers. Synthetic child speech samples are available on the project's GitHub. We also release two small synthetic child speech datasets generated from this work. Multi-speaker TTS proves to be a valuable approach for child speech synthesis, even with limited training data.

For Future work, we aim to perform a subjective evaluation (as proposed in [32]) on the released dataset for better clarity over the 'Naturalness', 'Intelligibility', and 'Speaker Similarity' of the generated child speech. Furthermore, it is also intended to investigate the use of synthetically generated child speech to enhance other areas of child speech research, such as ASR and speaker recognition.

ACKNOWLEDGMENT

The authors would like to acknowledge experts from Xperi: Gabriel Costache, Zoran Fejzo, Francisco Salgado, and George Sterpu for providing their expertise and feedback throughout. The authors would also like to thank Adriana Stan and Horia Cucu from the University Politehnica of Bucharest, for providing her expertise on TTS/ASR experiments.

REFERENCES

- [1] O. Watts, J. Yamagishi, K. Berkling, and S. King, 'HMM-based synthesis of child speech', *Proc 1st Workshop Child Computer Interaction ICMI08 Post-Conf. Workshop*, 2008.
- [2] O. Watts, J. Yamagishi, S. King, and K. Berkling, 'Synthesis of child speech with HMM adaptation and voice conversion', *IEEE Trans. Audio Speech Lang. Process.*, vol. 18, no. 5, pp. 1005–1016, 2010, doi: 10.1109/TASL.2009.2035029.
- [3] A. W. Black, H. Zen, and K. Tokuda, 'Statistical parametric speech synthesis', in *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, IEEE, 2007, p. IV-1229-IV-1232.
- [4] Maia, R., Zen, H., Gales, M.J.F. (2010) Statistical parametric speech synthesis with joint estimation of acoustic and excitation model parameters. *Proc. 7th ISCA Workshop on Speech Synthesis (SSW 7)*, 88-93
- [5] Y. Wang *et al.*, 'Tacotron: Towards End-To-End Speech Synthesis'.
- [6] J. Shen *et al.*, 'Natural TTS Synthesis By Conditioning Wavenet On Mel Spectrogram Predictions'.
- [7] Y. Ren *et al.*, 'Fastspeech: Fast, robust and controllable text to speech', *Adv. Neural Inf. Process. Syst.*, vol. 32, 2019.
- [8] Y. Ren *et al.*, 'FastSpeech 2: Fast and High-Quality End-to-End Text to Speech'. arXiv, Aug. 07, 2022.
- [9] C. Miao, S. Liang, M. Chen, J. Ma, S. Wang, and J. Xiao, 'Flow-TTS: A Non-Autoregressive Network for Text to Speech Based on Flow', in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 7209–7213. doi: 10.1109/ICASSP40776.2020.9054484.
- [10] J. Kim, S. Kim, J. Kong, and S. Yoon, 'Glow-TTS: A Generative Flow for Text-to-Speech via Monotonic Alignment Search'. Available: <https://github.com/jaywalnut310/glow-tts>.
- [11] A. Łańcucki, 'Fastpitch: Parallel text-to-speech with pitch prediction', in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2021, pp. 6588–6592.
- [12] M. Chen *et al.*, 'Adaspeech: Adaptive text to speech for custom voice', *ArXiv Prepr. ArXiv210300993*, 2021.
- [13] S. Arik, "Deep voice 2: multi-speaker neural text-to-speech," in *Proc. 31st Int. Conf. Neural Inf. Process. Syst.* Red Hook, NY, USA, Curran Associates Inc., 2017, pp. 2966–2974.
- [14] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, 'X-Vectors: Robust DNN Embeddings for Speaker Recognition', in *ICASSP, IEEE International Conference on Acoustics, Speech and*

Signal Processing - Proceedings, 2018, pp. 5329–5333. doi: 10.1109/ICASSP.2018.8461375.

- [15] L. Wan, Q. Wang, A. Papir, I. Lopez, and M. Moreno, ‘Generalized End-To-End Loss for Speaker Verification’ in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, Apr. 2018, pp. 4879–4883.
- [16] G. Heigold, I. Moreno, S. Bengio, and N. Shazeer, ‘‘End-to-end text-dependent speaker verification,’’ in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016, pp. 5115–5119.
- [17] M. Chen *et al.*, ‘Cross-lingual, Multi-speaker Text-To-Speech Synthesis Using Neural Speaker Embedding’, 2019, doi: 10.21437/Interspeech.2019-1632.
- [18] R. Valle, J. Li, R. Prenger, and B. Catanzaro, ‘‘Mellotron: Multispeaker expressive voice synthesis by conditioning on rhythm, pitch and global style tokens,’’ in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2020, pp. 6189–6193.
- [19] Y. Jia *et al.*, ‘Transfer learning from speaker verification to multispeaker text-to-speech synthesis’, in *Advances in Neural Information Processing Systems*, Neural information processing systems foundation, 2018, pp. 4480–4490.
- [20] A. Kulkarni, V. Colotte, and D. Jouvet, ‘‘Improving latent representation for end-to-end multispeaker expressive text to speech system,’’ Tech. Rep. fhal-02978485v1f, 2020.
- [21] E. Cooper, C.-I. Lai, Y. Yasuda, and J. Yamagishi, ‘‘Can speaker augmentation improve multi-speaker end-to-end TTS?’’ in *Proc. Interspeech*, Oct. 2020, pp. 1–5.
- [22] S. Shahnawazuddin, N. Adiga, and H. K. Kathania, ‘Effect of Prosody Modification on Children’s ASR’, *IEEE Signal Process. Lett.*, vol. 24, no. 11, pp. 1749–1753, Nov. 2017, doi: 10.1109/LSP.2017.2756347.
- [23] G. Yeung, R. Fan, and A. Alwan, ‘‘Fundamental frequency feature normalization and data augmentation for child speech recognition,’’ in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*, Toronto, ON, Canada, 2021, pp. 6993–6997, doi: 10.1109/ICASSP39728.2021.9413801.
- [24] S. Lee, A. Potamianos, and S. Narayanan, ‘‘Acoustics of children’s speech: Developmental changes of temporal and spectral parameters,’’ *J. Acoust. Soc. Amer.*, vol. 105, no. 3, pp. 1455–1468.
- [25] S. Shahnawazuddin, R. Sinha, and G. Pradhan, ‘Pitch-Normalized Acoustic Features for Robust Children’s Speech Recognition’, *IEEE Signal Process. Lett.*, vol. 24, no. 8, pp. 1128–1132, Aug. 2017, doi: 10.1109/LSP.2017.2705085.
- [26] S. Lee, A. Potamianos, and S. S. Narayanan, ‘Analysis of children’s speech: duration, pitch and formants’, in *EUROSPEECH*, 1997.
- [27] M. Gerosa, D. Giuliani, S. Narayanan, and A. Potamianos, ‘‘A review of ASR technologies for children’s speech,’’ in *Proc. 2nd Workshop Child, Comput. Interact. (WOCCI)*, 2009, pp. 1–8.
- [28] G. Stemmer, C. Hacker, S. Steidl, and E. Nöth, ‘‘Acoustic normalization of children’s speech,’’ in *Proc. 8th Eur. Conf. Speech Commun. Technol.*, 2003.
- [29] M. Gerosa, S. Lee, D. Giuliani, and S. Narayanan, ‘‘Analyzing children’s speech: An acoustic study of consonants and consonant-vowel transition,’’ in *Proc. IEEE Int. Conf. Acoustics Speech Signal Process.*, 2006, pp. I–I, doi: 10.1109/ICASSP.2006.1660040.
- [30] W. Ward, R. Cole, and S. Pradhan, ‘‘My science tutor and the MyST corpus,’’ Tech. Rep., 2019.
- [31] R. Jain, A. Barcovschi, M. Yiwere, D. Bigioi, P. Corcoran, and H. Cucu, ‘A Wav2vec2-Based Experimental Study on Self-Supervised Learning Methods to Improve Child Speech Recognition’, Apr. 2022, doi: 10.48550/arxiv.2204.05419.
- [32] R. Jain, M. Y. Yiwere, D. Bigioi, P. Corcoran, and H. Cucu, ‘A Text-to-Speech Pipeline, Evaluation Methodology, and Initial Fine-Tuning Results for Child Speech Synthesis’, *IEEE Access*, vol. 10, pp. 47628–47642, 2022, doi: 10.1109/ACCESS.2022.3170836.
- [33] K. Yang, T.-Y. Hu, J.-H. R. Chang, H. Swetha Koppula, and O. Tuzel, ‘Text is all You Need: Personalizing ASR Models Using Controllable Speech Synthesis’, in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Jun. 2023, pp. 1–5. doi: 10.1109/ICASSP49357.2023.10096971.
- [34] A. Fazel *et al.*, ‘SynthASR: Unlocking Synthetic Data for Speech Recognition’, in *Interspeech 2021*, ISCA, Aug. 2021, pp. 896–900. doi: 10.21437/Interspeech.2021-1882.
- [35] H. Zen *et al.*, ‘LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech’. arXiv, Apr. 05, 2019.
- [36] ‘IEEE Recommended Practice for Speech Quality Measurements’, *IEEE Trans. Audio Electroacoustics*, vol. 17, no. 3, pp. 225–246, 1969, doi: 10.1109/TAU.1969.1162058.
- [37] ‘The LJ Speech Dataset’. <https://keithito.com/LJ-Speech-Dataset>.
- [38] R. Badlani, A. \Lańcucki, K. J. Shih, R. Valle, W. Ping, and B. Catanzaro, ‘One TTS alignment to rule them all’, in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2022, pp. 6092–6096.
- [39] P. Neekhara, J. Li, and B. Ginsburg, ‘Adapting TTS models For New Speakers using Transfer Learning’. arXiv, Apr. 05, 2022. Available: <http://arxiv.org/abs/2110.05798>
- [40] F. Lux, J. Koch, and N. T. Vu, ‘Exact Prosody Cloning in Zero-Shot Multispeaker Text-to-Speech’, in *2022 IEEE Spoken Language Technology Workshop (SLT)*, Jan. 2023, pp. 962–969. doi: 10.1109/SLT54892.2023.10022433.
- [41] C.-P. Hsieh, S. Ghosh, and B. Ginsburg, ‘Adapter-Based Extension of Multi-Speaker Text-to-Speech Model for New Speakers’. arXiv, Nov. 01, 2022. Available: <http://arxiv.org/abs/2211.00585>
- [42] R. Prenger, R. Valle, and B. Catanzaro, ‘‘Waveglow: A flow-based generative network for speech synthesis,’’ in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2019, pp. 3617–3621.
- [43] W. Jang, D. Lim, and J. Yoon, ‘Universal melgan: A robust neural vocoder for high-fidelity waveform generation in multiple domains’, *ArXiv Prepr. ArXiv201109631*, 2020.
- [44] Y. Jiao, A. Gabryś, G. Tinchev, B. Putrycz, D. Korzekwa, and V. Klimkov, ‘Universal Neural Vocoding with Parallel Wavenet’, in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Jun. 2021, pp. 6044–6048. doi: 10.1109/ICASSP39728.2021.9414444.
- [45] J. Lorenzo-Trueba *et al.*, ‘Towards achieving robust universal neural vocoding’, *ArXiv Prepr. ArXiv181106292*, 2018.
- [46] Y. Jiao, A. Gabryś, G. Tinchev, B. Putrycz, D. Korzekwa, and V. Klimkov, ‘Universal neural vocoding with parallel wavenet’, in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2021, pp. 6044–6048.
- [47] D. Paul, Y. Pantazis, and Y. Stylianou, ‘Speaker Conditional WaveRNN: Towards Universal Neural Vocoder for Unseen Speaker and Recording Conditions’, Accessed: Mar. 10, 2022. [Online]. Available: <https://github.com/fatchord/WaveRNN>
- [48] P. L. Tobing and T. Toda, ‘High-fidelity and low-latency universal neural vocoder based on multiband WaveRNN with data-driven linear prediction for discrete waveform modeling’, 2021.
- [49] C.-C. Lo *et al.*, ‘MOSNet: Deep Learning based Objective Assessment for Voice Conversion’, in *Interspeech 2019*, Sep. 2019, pp. 1541–1545. doi: 10.21437/Interspeech.2019-2003.
- [50] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, ‘wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations’, in *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 2020, pp. 12449–12460.