

Deep Speech Synthesis from Articulatory Representations

Peter Wu¹, Shinji Watanabe², Louis Goldstein³, Alan W Black^{*2}, Gopala K. Anumanchipalli^{*1}

¹University of California, Berkeley, United States

²Carnegie Mellon University, United States

³University of Southern California, United States

peterw1@berkeley.edu

Abstract

In the articulatory synthesis task, speech is synthesized from input features containing information about the physical behavior of the human vocal tract. This task provides a promising direction for speech synthesis research, as the articulatory space is compact, smooth, and interpretable. Current works have highlighted the potential for deep learning models to perform articulatory synthesis. However, it remains unclear whether these models can achieve the efficiency and fidelity of the human speech production system. To help bridge this gap, we propose a time-domain articulatory synthesis methodology and demonstrate its efficacy with both electromagnetic articulography (EMA) and synthetic articulatory feature inputs. Our model is computationally efficient and achieves a transcription word error rate (WER) of 18.5% for the EMA-to-speech task, yielding an improvement of 11.6% compared to prior work. Through interpolation experiments, we also highlight the generalizability and interpretability of our approach.

Index Terms: speech synthesis, articulatory synthesis

1. Introduction

Speech synthesis has seen rapid development in recent years with deep learning based techniques. These models have shown success in text-to-speech (TTS) [1, 2, 3], speech-to-speech translation (S2ST) [4, 5, 6], voice conversion (VC) [7, 8, 9], and tasks with other modalities [10, 11, 12]. Moreover, this technology has yielded impactful technologies like speech synthesis aids for people with blindness or paralysis [13, 14, 10]. While speech synthesizers have already shown promising results in multiple domains, technologies like brain-to-speech devices remain challenging to build [10]. Thus, these unsolved tasks require new algorithms in order to achieve the development of high-fidelity, open-vocabulary synthesizers. To this end, our work focuses on devising a deep speech synthesis methodology that is computationally efficient, real-time, and high-fidelity. We propose a time-domain articulatory synthesis approach that is suitable for attaining these three properties and empirically validate our method on two distinct articulatory modalities, EMA and a synthetic articulatory space. Our deep learning models also exhibit valuable interpretability properties, which we demonstrate through interpolation experiments. Audio samples, code, and additional related information are available at <https://articulatorysynthesis.github.io>.

*Equal advising.

2. Speech Synthesis

2.1. Deep Speech Synthesis

Currently, state-of-the-art speech synthesis algorithms use deep learning [2, 10, 15, 7, 12]. While existing methods can generate high-fidelity speech, they tend to be computationally expensive and difficult to interpret and generalize [16, 17]. We attribute underspecification to the primary cause of these issues, as speech data is very high dimensional and current algorithms lack sufficient inductive biases. To help bridge this gap, we devise deep articulatory synthesis techniques that exhibit suitable computational efficiency, generalizability, and interpretability properties by behaving more similarly to the human speech production process than existing methods.

2.2. Articulatory Synthesis

Articulatory synthesis generally refers to the task of synthesizing speech from articulatory features, i.e., features containing information about the physical behavior of the human vocal tract [18, 19, 20, 21]. We identify two primary research directions in articulatory synthesis: 1. modelling the human vocal tract [22, 23, 24], and 2. learning the mapping from articulatory features to speech through a statistical means [25, 26, 27]. The former direction, due to its focus on computational modelling, has yielded articulatory synthesizers that are interpretable and relatively space-efficient but computationally slow. On the other hand, the latter direction has yielded methods that are much faster but have worse interpretability and memory efficiency. Ideally, speech synthesizers should have low space and time complexities, which would enable many impactful real-time applications. For example, such systems could allow patients with paralysis or aphasia to communicate naturally at any moment in time. Thus, we focus on making methods in the second research direction more memory-efficient in this work. Additionally, we highlight how statistical articulatory synthesis methods could also be highly interpretable, thus containing all of the benefits of articulatory synthesizers built using physical modelling.

Another motivation for our statistical research direction is the transferability of our methodology to all forms of speech synthesis. Current state-of-the-art speech synthesis systems rely on an intermediate speech representation, typically a spectrum or a learned representation [28, 29, 30, 31, 32]. Inductive biases offer one potential way of making these models efficient, generalizable, and interpretable, as mentioned in Section 2.1. Constraining these intermediate representations to an articulatory feature space is one way to impose such an inductive bias, since a limited set of articulator configurations can completely specify all possible human speech [33]. The resulting model would then need to perform an articulatory-to-speech mapping, of which the behavior is relatively unknown to our knowledge.

This work aims to bridge this gap by studying the efficiency, generalizability, interpretability, and fidelity of such a mapping using two distinct articulatory modalities. Specifically, we use the MNGU0 EMA dataset [34] and another corpus generated by VocalTractLab [24], detailed in Section 4.

While deep EMA-to-speech models have been previously studied, as far as we are aware [35, 36, 37], current models are not highly intelligible, achieving a transcription WER of around 30% on open-vocabulary tasks [35]. In this work, we build an EMA-to-speech model that achieves a transcription WER of 18.5% and perform detailed error analyses on the synthesized utterances. We also extend this approach to building a speech synthesizer using a synthetic articulatory modality. This model is efficient, high-fidelity, and interpretable, which has previously been unattained to our knowledge. We detail these models and our proposed time-domain articulatory synthesis methodology in Section 3 below.

3. Deep Articulatory Models

3.1. Frequency- and Time-Domain Modeling

Similarly to the state-of-the-art speech synthesis works discussed in Section 2, current deep articulatory synthesis works rely on synthesizing an intermediate spectrum representation, from which waveforms are generated [38, 39]. Since this behavior is not present in the human speech production process, we propose a model that directly maps articulatory representations to waveforms. We refer to this approach as a time-domain one, as it does not explicitly rely on a frequency-based intermediate. This modification noticeably improves model efficiency while achieving comparable intelligibility, as discussed in Sections 5 and 7. We detail our spectrum-intermediate baseline in Section 3.2 and our time-domain method in Section 3.3.

3.2. Spectrum-Intermediate Baseline

For our baseline, we build on a state-of-the-art model proposed by Gaddy and Klein [12]. Namely, we map articulatory representations to spectrums using a six-layer Transformer [40] prepended with three residual convolution blocks. To map spectrums to waveforms, we use HiFi-GAN [28], which we make autoregressive using the audio encoder from CAR-GAN [29]. For our spectrum representation, we use Mel spectrograms instead of MFCCs, as done in the CAR-GAN paper and most deep speech synthesis works [29, 1, 2]. We omit the phonemic loss to avoid requiring phoneme annotations during training and instead improve model performance by adding the HiFi-GAN adversarial loss [28]. Since articulatory representations in this work are pre-aligned with waveforms, we also omit the dynamic time warping loss. We refer to this model as the spectrum-intermediate (Spec.-Int.). Further modeling details can be found in the accompanying codebase.

3.3. HiFi-CAR Model

For our time-domain model, we feed our articulatory input features directly into HiFi-GAN [28], which we make autoregressive using the audio encoder from CAR-GAN [29]. To our knowledge, directly feeding articulatory inputs into a deep vocoder architecture has not yielded any successful results previously. However, we observe that this model is comparable to our baseline, as discussed in Section 7. Moreover, removing the need for an articulatory-to-spectrum architecture noticeably improves computational efficiency, as discussed in Section 5.

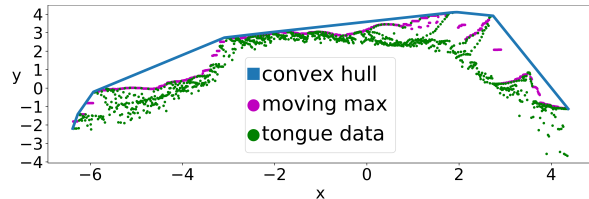


Figure 1: *Estimating the palate location with a convex hull and a moving maximum on tongue data in the training set (Sec. 4.1).*

We refer to this model as HiFi-CAR below. Further modeling details can be found in the accompanying codebase.

4. Datasets

4.1. Electromagnetic Articulography (EMA)

For our first task, we perform EMA-to-speech using the MNGU0 dataset [34]. MNGU0 contains 67 minutes of 16 kHz, single-speaker speech collected with the speaker instrumented in an EMA machine providing 200 Hz samples of EMA features. These 12-dimensional features contain the midsagittal x and y coordinates of jaw, lip, and tongue positions. We use the original train-test split with 1,129 training utterances and 60 test ones, and randomly choose 60 training datapoints for the validation set. Since EMA features do not contain voicing information, we concatenate them with estimated F0 sequences extracted using CREPE [41, 29].

Palate information, which is important for determining consonant sounds [42], is also not directly present in EMA. Thus, we estimate the location of the palate using tongue data in the training set using two methods. First, we compute the convex hull of the tongue coordinates [33]. Since this estimate does not account for the concave portions of the tongue, we also compute the moving maximum along the x-axis, using a window size of 32. For each EMA feature sequence, we estimate the tongue-to-palate distance by subtracting the tongue tip, body, and dorsum y-coordinates from both palate estimates, yielding 6 additional features. Figure 1 plots our tongue training data and palate estimates on MNGU0’s normalized xy-coordinates.

4.2. Synthetic Articulatory Features

Since EMA data does not explicitly contain enough manner information to perfectly reconstruct the original speech [10], we also experiment with synthetic articulatory data that does. Namely, we use the vocal tract model from Birkholz et al. [24] to create a single-speaker corpus of pseudo-words, each composed of two to three vowel and consonant sounds. Our training set has 10,000 such utterances, and our validation set has 250, totaling a few hours of speech. For our test set, we use vocal tract model outputs corresponding to the first 99 phoneme sequences in the CMU US KAL Diphone database [43]. All waveforms have a sampling rate of 44100 Hz and 30-dimensional articulatory features are recorded every 110 samples. We refer to this dataset as the Birkholz-Pseudoword (Birk.-Pseudo.) dataset below.

5. Computational Efficiency

Computational efficiency during inference is essential for building real-time, on-device speech synthesizers. We observe that our time-domain articulatory synthesis model is more time- and

Data	CPU (s) ↓	GPU (s) ↓	Params. ↓
HiFi-CAR	4.10 ± 0.03	0.37 ± 0.00	1.3 * 10 ⁷
Spec.-Int.	7.61 ± 0.02	0.39 ± 0.00	9.5 * 10 ⁷

Table 1: Average inference time and number of parameters for EMA-to-speech models. See Section 5 for details.

space-efficient than the spectrum-intermediate baseline. Table 1 contains the EMA-to-speech inference times and number of parameters for both models. GPU trials use one RTX 2080 Ti GPU, and CPU trials use none. We report inference time as the mean and standard deviation of five trials, each calculating the average time to synthesize an utterance in our test set. Our time-domain model is almost $2\times$ faster than the baseline in the CPU-only case and uses $7\times$ less parameters. Our synthetic articulatory experiments, which use similar hyperparameters, yield matching trends, detailed in the supplementary material linked in Section 1. These results suggest that directly mapping articulatory features to speech is more efficient than relying on an intermediate spectral representation.

6. Interpolation

6.1. Vowel Interpolation

To study the generalizability of our time-domain model, we perform interpolation experiments, prompting our model to synthesize unseen articulatory representations between pairs of sounds. All generations are available through our accompanying link. First, we interpolate between “ta”, “tu”, and “ti” to analyze how well our model generalizes across vowels. Similarly to our synthetic articulatory task, we generate the articulatory features for these sounds using VocalTractLab [24], which can provide versions of these utterances with the same duration. For each of the three possible pairs of sounds, we perform a linear interpolation between the two articulatory features, generating seven evenly spaced weighted combinations. Figure 2 contains the mel-spectrograms of the generated speech from our model for each of these combined articulatory features. The transitions between spectrum values in each interpolation are smooth, suggesting that our network is able to model the continuity of articulator movements, at least with respect to vowels. This trend is also reflected in our listening test, where English-speaking, MTurk listeners classified each of the seven “ta”→“tu” utterances as “ta”, “te”, “tu”, or “possibly” followed by one of these three syllables. Each utterance received 10 votes, yielding 70 annotations by a total of 19 listeners. To calculate which syllable the listeners assigned to each utterance on average, we convert each vote to a size-3 one-hot vector, multiplied by 0.5 if “possibly” was used. Index 0 is non-zero if “ta” was chosen, 1 if “te”, and 2 if “tu”. For each utterance, we plot the mean vector as three adjacent bars, one for each scalar, depicted in Figure 4. In this transition from “ta” to “tu”, listeners steadily perceive “ta” less and “tu” more, with “te” peaking in the middle, matching human vowel speech production behavior [44].

6.2. Consonant Interpolation

We also study the generalizability of our model with respect to consonants. To study how well our model generalizes across types of consonant sounds, we fix the place of articulation and interpolate between consonant types. Namely, we interpolate between the alveolar consonants “ra”, “na”, and “la”, using the same generation and evaluation methodologies as

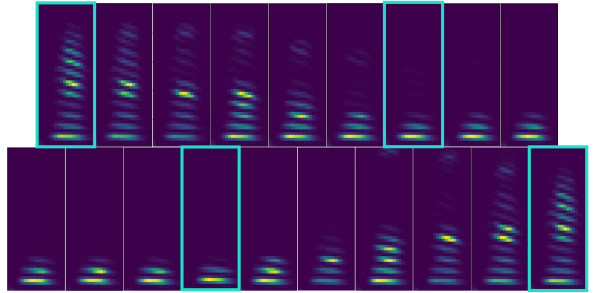


Figure 2: Vowel interpolation between “ta” (top left), “tu”, “ti”, and “ta” (bottom right), with these four diphones boxed.

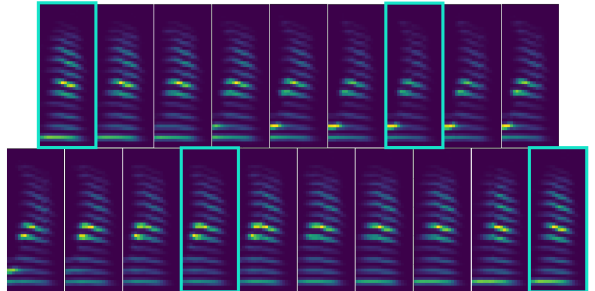


Figure 3: Alveolar consonant interpolation between “ra” (top left), “na”, “la”, and “ra” (bottom right), with these four diphones boxed.

our vowel interpolation experiment. Figure 3 depicts the mel-spectrograms of the synthesized interpolation samples. Our time-domain model smoothly transitions between these nasal, approximant, and lateral approximant consonants, as done in the human speech production process. MTurk listeners also perceived this behavior, based on Figure 4 results similar to those in Section 6.1. Thus, while sounds between these diphones are not in our training set, our model is still able to generate unseen transitions that reflect those made by human articulators, suggesting an ability to generalize.

6.3. Interpretability

We note that these interpolation results also highlight the interpretability of articulatory features. Namely, we are able to simply take an element-wise weighted sum of two same-length sequences of articulatory features in order to create the utterance corresponding to articulator movements in between the two gestures. For example, to create the “te” sound, we would just need to synthesize the average of the articulatory feature sequences for “ti” and “ta”. To our knowledge, this degree of interpretability is not supported by other speech representations like spectrums or deep-learning-based ones [45, 46].

7. Synthesis Quality

7.1. Subjective Fidelity Evaluation

We first compare the fidelity our models by synthesizing the EMA-to-speech test set utterances and performing an AB naturalness preference test on MTurk with 40 total English-speaking listeners and 2 comparisons per datapoint, yielding 120 votes. Samples are available in the supplementary material linked in Section 1. Votes were evenly split between the models, with

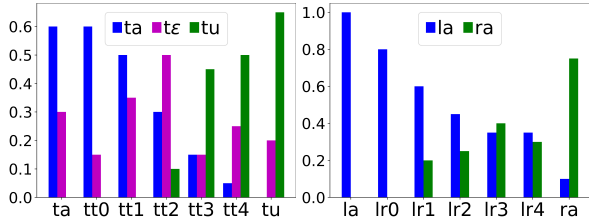


Figure 4: Human perception of “ta”→“tu” (left) and “la”→“ra” (right) interpolations. Details in Section 6.

Model	MCD ↓	
	Birk.-Pseudo	EMA-MGNU0
HiFi-CAR	3.26 ± 0.18	4.81 ± 0.73
Spec.-Int.	5.15 ± 0.48	4.75 ± 0.81

Table 2: MCD for each model on Birkholz and EMA data. See Section 7.2 for details.

listeners agreeing on 28 out of 60 pairs, suggesting that listeners are unable to distinguish the two models. Thus, to study the differences between our models, we perform transcription and objective fidelity evaluations, as discussed below.

7.2. Objective Fidelity Evaluation

In order to perform an objective evaluation of synthesis quality, we compute the mel-cepstral distortion (MCD) [47] between the ground truth and the generations from both models for each utterance in our test sets. As described in Table 2, our time-domain articulatory synthesis approach performs better than the spectrum-intermediate baseline on the synthetic articulatory dataset and slightly worse on the EMA-to-speech task. We attribute the performance drop of our model on the EMA task to information loss within in the input data. Namely, the model appears to confuse phonemes due to the lack of manner information in the EMA inputs, as heard in the accompanying samples and analyzed in Section 8.

7.3. Transcription

To evaluate intelligibility, we conduct open-vocabulary transcription experiments using the EMA-to-speech task, for which the test set contains non-synthetic utterances. First, we perform an objective automatic speech recognition (ASR) evaluation using DeepSpeech [48] and Wav2Vec2 [49, 50]. We use ASR to transcribe the synthesis outputs of our models on the entire MNGU0 evaluation set described in Section 4.1 and calculate the average word error rates (WERs). We also evaluate the intelligibility of our model through human evaluations, following the MTurk methodology in Taguchi and Kaburagi’s EMA-to-speech work [35] in order to compare with their WER 30.1% result. Table 3 summarizes our WERs. Based on 120 transcriptions per model by a total of 70 English-speaking listeners, our models achieve an average WER between 18% and 19%, i.e., improving WER by over 11%. Our models have comparable WERs for both the objective and human evaluations, suggesting that they are both suitable for generating intelligible speech.

8. Phoneme Confusion

To study the phonological errors made by our model, we analyze the phonemes that our time-domain EMA-to-speech model

Model	WER ↓		
	Human	DeepSpeech	Wav2Vec2
HiFi-CAR	18.5	34.7	21.7
Spec.-Int.	18.9	31.0	23.5

Table 3: WER for each model. See Section 7.3 for details.

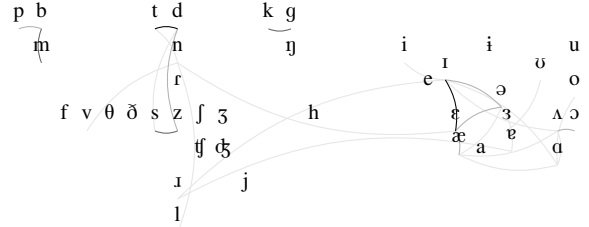


Figure 5: Phoneme confusability based on ASR transcriptions. Phoneme pairs that are confused more frequently have darker lines. Details in Section 8.

confused during synthesis. For all ASR-transcribed test utterances and their ground truth texts, we convert the graphemes to phoneme sequences using Phonemizer [51]. We identify phoneme confusion pairs using *scLite*,¹ which aligns each predicted sequence with the respective ground truth and then records the substitution errors. Figure 5 plots pairs that are confused more than once on an International Phonetic Alphabet (IPA) chart [12], where pairs with higher frequencies of substitution errors have darker lines. Most of the word substitution errors are due to plosive or vowel confusions. One potential reason for the plosive substitutions is that plosives generally have a shorter duration than other consonant types like fricatives [52] and thus may be more readily confusable. Multiple voiced-unvoiced pairs are also confused, which may be because estimated F0 is the only voicing information in the input, as described in Section 4.1. These results suggest that future work adding additional articulatory features like velar and glottal information can yield even higher fidelity articulatory synthesizers, potentially comparable to current text-to-speech (TTS) systems.

9. Conclusion and Future Directions

In this work, we study ways to build deep articulatory synthesizers that are efficient and high-fidelity. Based on computational efficiency evaluations, we observe that our proposed time-domain methodology is suitable for achieving time and space complexities that are noticeably lower than the baseline spectrum-intermediate approach. Our interpolation study also highlights the generalizability and interpretability of our approach. Through MCD, ASR, and human transcription experiments, we demonstrate the intelligibility of our model, improving the transcription word error rate for the EMA-to-speech task by over 11%. Moving forward, we plan to test our methodology on other modalities like electromyography [12] and real-time magnetic resonance imaging [53].

10. Acknowledgements

This research is supported by the National Science Foundation (Award 2106928).

¹<https://github.com/usnistgov/SCTK>

11. References

- [1] Y. Wang *et al.*, “Tacotron: Towards end-to-end speech synthesis,” in *Interspeech*, 2017.
- [2] T. Hayashi *et al.*, “Espnet2-tts: Extending the edge of tts research,” *arXiv preprint arXiv:2110.07840*, 2021.
- [3] R. Prenger, R. Valle, and B. Catanzaro, “Waveglow: A flow-based generative network for speech synthesis,” in *ICASSP*, 2019.
- [4] A. Tjandra, S. Sakti, and S. Nakamura, “Speech-to-speech translation between untranscribed unknown languages,” in *ASRU*, 2019, pp. 593–600.
- [5] Y. Jia *et al.*, “Direct speech-to-speech translation with a sequence-to-sequence model,” in *Interspeech*, 2019.
- [6] H. Inaguma *et al.*, “ESPnet-ST: All-in-one speech translation toolkit,” in *ACL*, 2020.
- [7] A. Polyak *et al.*, “Speech Resynthesis from Discrete Disentangled Self-Supervised Representations,” in *Interspeech*, 2021.
- [8] P. Wu *et al.*, “Understanding the tradeoffs in client-side privacy for downstream speech tasks,” in *APSIPA ASC*, 2021.
- [9] B. Sisman *et al.*, “An overview of voice conversion and its challenges: From statistical modeling to deep learning,” *TASLP*, 2020.
- [10] G. K. Anumanchipalli, J. Chartier, and E. F. Chang, “Speech synthesis from neural decoding of spoken sentences,” *Nature*, vol. 568, no. 7753, pp. 493–498, 2019.
- [11] C. Yu *et al.*, “Dorian: Duration informed attention network for multimodal synthesis,” *arXiv preprint arXiv:1909.01700*, 2019.
- [12] D. Gaddy and D. Klein, “An improved model for voicing silent speech,” in *ACL-IJCNLP*, 2021.
- [13] A. Karmel *et al.*, “IoT based assistive device for deaf, dumb and blind people,” *ICRTAC*, 2019.
- [14] M. Angrick *et al.*, “Speech synthesis from ecog using densely connected 3d convolutional neural networks,” *Journal of neural engineering*, vol. 16, no. 3, p. 036019, 2019.
- [15] Y. Jia *et al.*, “Translatotron 2: Robust direct speech-to-speech translation,” *arXiv preprint arXiv:2107.08661*, 2021.
- [16] T. Nekvinda and O. Dušek, “One Model, Many Languages: Meta-Learning for Multilingual Text-to-Speech,” in *Interspeech*, 2020.
- [17] Y. Zhang *et al.*, “Learning to speak fluently in a foreign language: Multilingual speech synthesis and cross-language voice cloning,” *arXiv preprint arXiv:1907.04448*, 2019.
- [18] G. Fant, “What can basic research contribute to speech synthesis?” *Journal of Phonetics*, vol. 19, no. 1, pp. 75–90, 1991.
- [19] P. Rubin, T. Baer, and P. Mermelstein, “An articulatory synthesizer for perceptual research,” *The Journal of the Acoustical Society of America*, vol. 70, no. 2, pp. 321–328, 1981.
- [20] C. Scully, “Articulatory synthesis,” in *Speech production and speech modelling*. Springer, 1990, pp. 151–186.
- [21] J. Lian, A. W. Black, L. Goldstein, and G. K. Anumanchipalli, “Deep neural convolutive matrix factorization for articulatory representation decomposition,” *Interspeech*, 2022.
- [22] G. Fant, “The lf-model revisited. transformations and frequency domain analysis,” *Speech Trans. Lab. Q. Rep., Royal Inst. of Tech. Stockholm*, vol. 2, no. 3, p. 40, 1995.
- [23] K. Iskarous *et al.*, “Casy: The haskins configurable articulatory synthesizer,” in *International Congress of Phonetic Sciences*, 2003.
- [24] P. Birkholz, “Modeling consonant-vowel coarticulation for articulatory speech synthesis,” *PLoS one*, vol. 8, no. 4, p. e60603, 2013.
- [25] S. Aryal and R. Gutierrez-Osuna, “Data driven articulatory synthesis with deep neural networks,” *Computer Speech & Language*, vol. 36, pp. 260–273, 2016.
- [26] F. Bocquelet *et al.*, “Robust articulatory speech synthesis using deep neural networks for bci applications,” in *Interspeech*, 2014.
- [27] Y.-W. Chen *et al.*, “Ema2s: An end-to-end multimodal articulatory-to-speech system,” in *ISCAS*, 2021.
- [28] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” in *NeurIPS*, 2020.
- [29] M. Morrison *et al.*, “Chunked autoregressive gan for conditional waveform synthesis,” in *Submitted to ICLR 2022*, April 2022.
- [30] R. Badlani *et al.*, “One tts alignment to rule them all,” *arXiv preprint arXiv:2108.10447*, 2021.
- [31] J. Kim, J. Kong, and J. Son, “Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech,” in *ICML*, 2021.
- [32] I. Elias *et al.*, “Parallel tacotron 2: A non-autoregressive neural tts model with differentiable duration modeling,” *ArXiv*, vol. abs/2103.14574, 2021.
- [33] H. Nam *et al.*, “A procedure for estimating gestural scores from speech acoustics,” *The Journal of the Acoustical Society of America*, vol. 132, no. 6, pp. 3980–3989, 2012.
- [34] K. Richmond, P. Hoole, and S. King, “Announcing the electromagnetic articulography (day 1) subset of the mngu0 articulatory corpus,” in *Interspeech*, 08 2011, pp. 1505–1508.
- [35] F. Taguchi and T. Kaburagi, “Articulatory-to-speech conversion using bi-directional long short-term memory,” in *Interspeech*, 2018, pp. 2499–2503.
- [36] S. Stone, P. Schmidt, and P. Birkholz, “Prediction of voicing and the f0 contour from electromagnetic articulography data for articulation-to-speech synthesis,” in *ICASSP*, 2020.
- [37] Z.-C. Liu, Z.-H. Ling, and L.-R. Dai, “Articulatory-to-acoustic conversion using blstm-rnns with augmented input representation,” *Speech Communication*, vol. 99, pp. 161–172, 2018.
- [38] T. G. Csap’o *et al.*, “Ultrasound-based articulatory-to-acoustic mapping with waveglow speech synthesis,” in *Interspeech*, 2020.
- [39] M.-A. Georges *et al.*, “Towards an articulatory-driven neural vocoder for speech synthesis,” in *International Seminar on Speech Production*, 2020.
- [40] A. Vaswani *et al.*, “Attention is all you need,” in *NeurIPS*, 2017.
- [41] J. W. Kim *et al.*, “Crepe: A convolutional representation for pitch estimation,” in *ICASSP*, 2018.
- [42] C. Hagedorn, M. Proctor, and L. Goldstein, “Automatic analysis of singleton and geminate consonant articulation using real-time magnetic resonance imaging,” in *Interspeech*, 2011.
- [43] K. Lenzo and A. Black, “Diphone collection and synthesis,” *ICSLP*, 2000.
- [44] P. J. Alfonso and T. Baer, “Dynamics of vowel articulation,” *Language and Speech*, vol. 25, no. 2, pp. 151–173, 1982.
- [45] Y. Wang *et al.*, “Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis,” in *ICML*, 2018.
- [46] G. Sun *et al.*, “Fully-hierarchical fine-grained prosody modeling for interpretable speech synthesis,” in *ICASSP*, 2020.
- [47] A. W. Black, “CMU wilderness multilingual speech dataset,” in *ICASSP*, 2019.
- [48] A. Y. Hannun *et al.*, “Deep speech: Scaling up end-to-end speech recognition,” *ArXiv*, vol. abs/1412.5567, 2014.
- [49] A. Baevski *et al.*, “Wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *NeurIPS*, 2020.
- [50] M. Ravanelli *et al.*, “SpeechBrain: A general-purpose speech toolkit,” 2021, arXiv:2106.04624.
- [51] M. Bernard and H. Titeux, “Phonemizer: Text to phones transcription for multiple languages in python,” *Journal of Open Source Software*, vol. 6, no. 68, p. 3958, 2021.
- [52] A. Alwan, J. Jiang, and W. Chen, “Perception of place of articulation for plosives and fricatives in noise,” *Speech communication*, vol. 53, no. 2, pp. 195–209, 2011.
- [53] Y. Lim *et al.*, “A multispeaker dataset of raw and reconstructed speech production real-time mri video and 3d volumetric images,” *Scientific Data*, vol. 8, 07 2021.