

# Deep Speech Synthesis from Articulatory Features

Anonymous ACL submission

## 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 both computationally efficient and highly intelligible, achieving a transcription word error rate (WER) of 7.14% for the EMA-to-speech task. Through interpolation experiments, we also highlight the generalizability and interpretability of our approach.

## 1 Introduction

Speech synthesis has seen rapid development in recent years with deep learning based techniques. These models have shown success in tasks like text-to-speech (TTS) (Wang et al., 2017; Hayashi et al., 2021; Prenger et al., 2019), speech-to-speech translation (S2ST) (Tjandra et al., 2019; Jia et al., 2019; Inaguma et al., 2020), voice conversion (VC) (Polyak et al., 2021; Wu et al., 2021a; Sisman et al., 2020), and more (Anumanchipalli et al., 2019; Yu et al., 2019; Gaddy and Klein, 2021). Moreover, this technology has yielded impactful technologies like speech synthesis aids for people with blindness or paralysis (Karmel et al., 2019; Angrick et al., 2019; Anumanchipalli et al., 2019). While speech synthesizers have already shown promising results for assistive tasks in healthcare and other challenging domains, technologies like brain-to-speech devices are still nascent and require new

algorithms in order to be deployed as 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.

We proceed by discussing speech synthesis in the context of deep learning and articulatory synthesis in Section 2. In Section 3, we describe our deep articulatory models and time-domain methodology. Then, we discuss the two articulatory datasets chosen for our empirical studies and their respective modalities in Section 4. With these datasets, we conduct computational efficiency, interpolation, and synthesis quality studies, discussed in Sections 5, 6, and 7, respectively. We then provide further analyses with respect to phoneme confusability in Section 8. Finally, we summarize our results and propose future directions in Section 9. Audio samples and additional related information are all available at <https://articulatorysynthesis.github.io>.

## 2 Speech Synthesis

### 2.1 Deep Speech Synthesis

Currently, state-of-the-art speech synthesis algorithms use deep learning (Hayashi et al., 2021; Anumanchipalli et al., 2019; Jia et al., 2021; Polyak et al., 2021; Gaddy and Klein, 2021). While existing methods can generate high-fidelity speech, they tend to be computationally expensive and difficult to interpret and generalize (Nekvinda and Dušek, 2020; Zhang et al., 2019). We attribute underspecification to the primary cause of these issues, as speech data is very high dimensional

081 and current algorithms lack sufficient inductive  
082 biases. To help bridge this gap, we devise deep  
083 articulatory synthesis techniques that exhibit suit-  
084 able computational efficiency, generalizability, and  
085 interpretability properties by behaving more simi-  
086 larly to the human speech production process than  
087 existing methods.

## 088 2.2 Articulatory Synthesis

089 Articulatory synthesis generally refers to the task  
090 of synthesizing speech from articulatory features,  
091 i.e., features containing information about the phys-  
092 ical behavior of the human vocal tract (Fant, 1991;  
093 Rubin et al., 1981; Scully, 1990). We identify two  
094 primary research directions in articulatory synthe-  
095 sis: 1. modelling the human vocal tract (Fant,  
096 1995; Iskarous et al., 2003; Birkholz, 2013a), and  
097 2. learning the mapping from articulatory fea-  
098 tures to speech through a statistical means (Aryal  
099 and Gutierrez-Osuna, 2016; Bocquelet et al., 2014;  
100 Chen et al., 2021). The former direction, due to  
101 its focus on computational modelling, has yielded  
102 articulatory synthesizers that are interpretable and  
103 relatively space-efficient but computationally slow.  
104 On the other hand, the latter direction has yielded  
105 methods that are much faster but have worse inter-  
106 pretability and memory efficiency. Ideally, speech  
107 synthesizers should have low space and time com-  
108 plexities, which would enable many impactful real-  
109 time applications. For example, such systems could  
110 allow patients with paralysis or aphasia to commu-  
111 nicate naturally at any moment in time. Thus, we  
112 focus on making methods in the second research  
113 direction more memory-efficient in this work. Ad-  
114 ditionally, we highlight how statistical articula-  
115 tory synthesis methods could also be highly inter-  
116 pretable, thus containing all of the benefits of  
117 articulatory synthesizers built using physical mod-  
118 elling.

119 We also focus on the statistical research direc-  
120 tion in this work because of the transferability of  
121 our methodology to all forms of speech synthesis.  
122 Current state-of-the-art speech synthesis systems  
123 rely on an intermediate speech representation, typi-  
124 cally a spectrum or a learned representation (Kong  
125 et al., 2020; Morrison et al., 2022; Badlani et al.,  
126 2021; Kim et al., 2021; Elias et al., 2021). Induc-  
127 tive biases offer one potential way of making these  
128 models efficient, generalizable, and interpretable as  
129 mentioned in Section 2.1. Constraining these inter-  
130 mediate representations to an articulatory feature

space is one way to impose such an inductive bias,  
especially since there is a limited set of articulator  
configurations that can completely specify all possi-  
ble human speech. The resulting model would  
then need to perform an articulatory-to-speech map-  
ping, of which the behavior is relatively unknown  
to our knowledge. This work aims to bridge this  
gap by studying the efficiency, generalizability, in-  
terpretability, and fidelity of such a mapping using  
two distinct articulatory modalities, EMA and a  
synthetic one generated using a vocal tract model,  
detailed in Section 4.

While deep EMA-to-speech models have been  
previously studied, as far as we are aware (Taguchi  
and Kaburagi, 2018; Stone et al., 2020; Liu et al.,  
2018), current models are not highly intelligible,  
achieving a transcription WER of around 30%  
on open-vocabulary tasks (Taguchi and Kaburagi,  
2018). In this work, we build an EMA-to-speech  
model that achieves a transcription WER of 7.14%  
and perform detailed error analyses on the synthe-  
sized 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 articula-  
tory synthesis works rely on synthesizing an inter-  
mediate spectrum representation, from which wave-  
forms are generated (Csap'o et al., 2020; Georges  
et al., 2020). Since this behavior is not present in  
the human speech production process, we propose  
a model that directly maps articulatory features to  
waveforms in this work. Since this model does  
not explicitly rely on a frequency-based interme-  
diate, we refer to this approach as a time-domain  
one. This modification noticeably improves model  
efficiency while achieving comparable intelligibil-  
ity on our two datasets, as discussed in Sections  
5 and 7. We proceed to discuss our spectrum-  
intermediate baseline in Section 3.2 and our two  
time-domain methods in Sections 3.3 and 3.4.

178	<b>3.2 Spectrum-Intermediate Baseline</b>	
179	For our baseline deep learning model, we build	228
180	on the state-of-the-art articulatory synthesis archi-	229
181	tecture proposed by Gaddy and Klein (Gaddy and	230
182	Klein, 2021). Namely, we map articulatory fea-	231
183	tures to spectrums using a six-layer Transformer	232
184	(Vaswani et al., 2017) prepended with three resid-	
185	ual convolution blocks. To map spectrums to	
186	waveforms, we use HiFi-GAN (Kong et al., 2020),	
187	which has been shown to perform better than the	
188	WaveNet vocoder used by Gaddy and Klein (Gaddy	
189	and Klein, 2021). For our spectrum representation,	
190	we use Mel spectrograms instead of MFCCs, as	
191	done in the HiFi-GAN paper and most deep speech	
192	synthesis works (Kong et al., 2020; Wang et al.,	
193	2017; Hayashi et al., 2021).	
194	We also modify the loss function used by Gaddy	
195	and Klein (Gaddy and Klein, 2021). To avoid re-	
196	quiring phoneme annotations to train the model, we	
197	omit the phonemic loss. We instead improve model	
198	performance by adding the adversarial loss used	
199	by HiFi-GAN (Kong et al., 2020). Since our data	
200	in this work has sequences of articulatory features	
201	that are pre-aligned with waveforms, we also do not	
202	need the dynamic time warping loss. We refer to	
203	this resulting baseline as the spectrum-intermediate	
204	(Spec.-Int.) model below.	
205	In all of our experiments, we train the Trans-	
206	former model using the Adam optimizer (Kingma	
207	and Ba, 2015) with a learning rate of $3.0 * 10^{-5}$ for	
208	both the generator and the discriminators, a batch	
209	size of 32, and loss balancing coefficients match-	
210	ing those used with the original HiFi-GAN model	
211	(Kong et al., 2020). Our discriminator architec-	
212	tures and HiFi-GAN spectrum-to-speech vocoder	
213	parameters also match those of Kong et al. (Kong	
214	et al., 2020), and our Transformer has a hidden	
215	dimension of 1024 and a dropout rate of 0.2.	
216	<b>3.3 Time-Domain HiFi-GAN</b>	
217	For our first time-domain model, we feed our ar-	
218	ticulatory input features directly into HiFi-GAN	
219	(Kong et al., 2020), keeping the architecture and	
220	loss functions the same while changing the input	
221	modality. To our knowledge, directly feeding ar-	
222	ticulatory inputs into a deep vocoder architecture	
223	has not yielded any successful results previously.	
224	However, we observe that this model is compa-	
225	rable to our baseline, as discussed in Section 7.	
226	Moreover, removing the need for an articulatory-	
227	to-spectrum architecture noticeably improves com-	
	putational efficiency, as discussed in Section 5. For	228
	all of our experiments, we optimize this model us-	229
	ing the same hyperparameters as the HiFi-GAN	230
	spectrum-to-speech vocoder used in the Section	231
	3.2 baseline above.	232
	<b>3.4 NSF-CAR Model</b>	233
	For our second time-domain model, we build on the	234
	neural source-filter (NSF) architecture (Wang et al.,	235
	2019). Since articulatory features can be divided	236
	into source- and filter-related attributes (Birkholz,	237
	2013a), we experiment with this architecture in	238
	order to study whether explicitly modelling this	239
	separation could improve articulatory synthesis per-	240
	formance.	241
	Similarly to our baseline, we use the loss func-	242
	tion from HiFi-GAN to improve synthesis fidelity.	243
	We also leverage autoregression to improve the	244
	pitch and periodicity of model outputs and make	245
	our model a streaming-based one. Namely, we	246
	incorporate the autoregressive encoder from CAR-	247
	GAN (Morrison et al., 2022) into our model, con-	248
	catenating its output with each vector in the condi-	249
	tion module input sequence. We replace the convo-	250
	lutions in the NSF condition module with GBlock	251
	layers (Morrison et al., 2022), which we found to	252
	further improve model performance. Figure 7 in	253
	the Appendix depicts the architecture of our gener-	254
	ator.	255
	To our knowledge, neural source filter mod-	256
	els are currently only used for building vocoders	257
	that map spectrums to speech (Wang et al., 2019;	258
	Georges et al., 2020). In this work, we leverage	259
	source-filter modelling to perform articulatory syn-	260
	thesis without relying on an intermediate spectrum	261
	representation.	262
	<b>3.5 WSOLA</b>	263
	As observed by Morrison et al. (Morrison et al.,	264
	2022), simply concatenating the output chunks gen-	265
	erated through an autoregressive process yields arti-	266
	facts at the concatenation points. Thus, during eval-	267
	uation, we join outputs using an approach based	268
	on WSOLA. Namely, we overlap-and-add adjacent	269
	output chunks at intersections with maximum cross-	270
	correlation, sliding the chunks up to a distance of	271
	one pitch period. We calculate a pitch period by	272
	multiplying the sampling rate with the reciprocal	273
	of the last F0 value in the first chunk input. Figure	274
	1 depicts one such WSOLA operation.	275

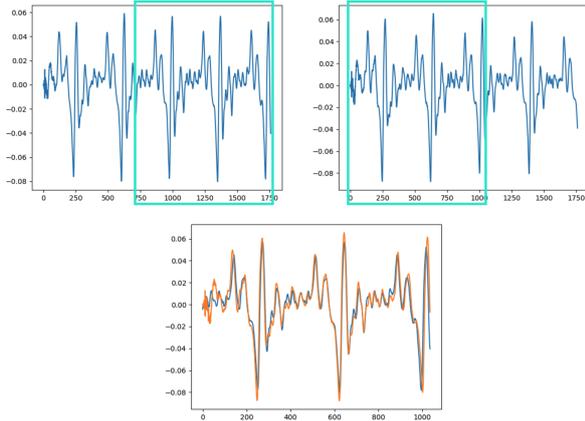


Figure 1: WSOLA-based method for concatenating waveforms.

## 4 Datasets

### 4.1 Electromagnetic Articulography (EMA)

For our first task, we perform EMA-to-speech using the MNGU0 dataset (Richmond et al., 2011), which contains 67 minutes of single-speaker speech recorded at 16 kHz annotated with 12-dimensional EMA features recorded at 200 Hz. We use the train-test split provided in the original work, which has 1,129 utterances for training and 60 for testing. Among the 1,129 training utterances, we set off a random size-60 subset for validation. Since EMA on its own does not contain voicing information, we concatenate estimated F0 sequences extracted using CREPE (Kim et al., 2018; Morrison et al., 2022) to the EMA features, forming a 13-dimensional input feature.

### 4.2 Synthetic Articulatory Features

Since EMA data does not contain enough manner information to perfectly reconstruct the original speech, we also experiment with synthetic articulatory data that does. Namely, we use the vocal tract model from Birkholz et al. (Birkholz, 2013a) 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 evaluation set, we use the Birkholz vocal tract model outputs corresponding to the first 99 phoneme sequences in the CMU US KAL Diphone database (Lenzo and Black, 2000). All waveforms have a sampling rate of 44100 Hz and articulatory features are recorded every 110 samples. We refer to this dataset as the Birkholz-Pseudoword (Birk.-Pseudo.) dataset be-

low. In this dataset, our articulatory features are 30-dimensional.

## 5 Computational Efficiency

Computational efficiency during training is essential for low-resource speech synthesis tasks like brain-to-speech and other articulatory synthesis tasks where data collection is expensive. During inference, computational efficiency is essential for building real-time speech synthesizers, e.g., for brain-to-speech. We observe that our time-domain articulatory synthesis model has some suitable computational efficiency properties compared to the frequency-domain baseline. As shown in Table 1, our model is able to train twice as fast as the baseline on a single RTX 2080 Ti GPU for the task with synthetic articulatory data. While our model synthesizes utterances slower than the baseline due to the nature of autoregression (Morrison et al., 2022), we observe that generation on a CPU is still faster than real-time.

Compared to the baseline, our time-domain models are much more memory efficient, as detailed in Table 2. Our models are able to use over 8 to 20 times less number of parameters than the baseline due to their ability to directly map articulatory features to speech. Namely, while current articulatory synthesis models like our baseline rely on two components, one to output spectrums and another to convert spectrums to waveforms, our time-domain models only contain one. We note that the real-time and memory efficient properties of our time-domain models make them a viable choice for streaming, on-device tasks.

Data	Birk.-Pseudo.	EMA-MGNU0
NSF-CAR	34	81
HiFi-GAN	8	9
Spec.-Int.	68	80

Table 1: Total training time for each model in hours.

Model	Birk.-Pseudo.	EMA-MGNU0
NSF-CAR	$4.4 * 10^6$	$4.2 * 10^6$
HiFi-GAN	$14.2 * 10^6$	$12.6 * 10^6$
Spec.-Int.	$98.7 * 10^6$	$94.0 * 10^6$

Table 2: Number of parameters of each model.

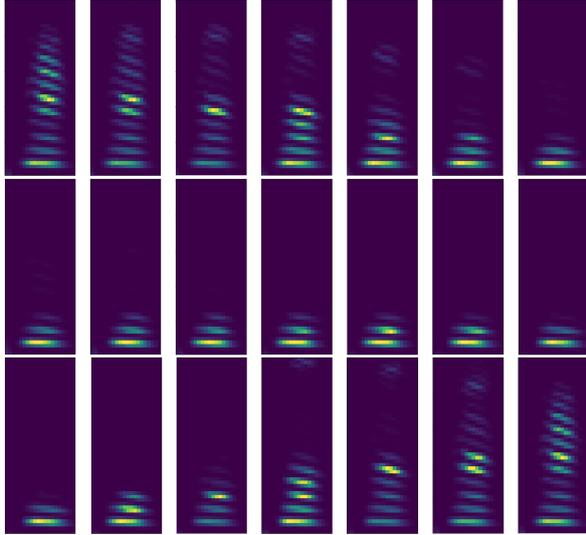


Figure 2: Vowel interpolation. The top row contains the synthesized samples between the "ta" and "tu" sounds, the middle row "tu" and "ti", and the bottom row "ti" and "ta".

## 6 Interpolation

### 6.1 Vowel Interpolation

To study the generalizability of our time-domain model, we perform interpolation experiments. First, to analyze how well our model generalizes across vowel sounds, we perform vowel interpolation. Namely, we interpolate between the "ta" and "tu" sounds, "tu" and "ti", and "ti" and "ta" using the synthetic articulatory data. We generate the articulatory features for "ta", "tu", and "ti" using the code provided by Birkholz et al., similarly to our approach for creating the synthetic articulatory dataset described above. For each of the three pairs of sounds, we perform a linear interpolation between the two articulatory features, generating seven evenly spaced weighted combinations. The figures below are generated using outputs from our NSF-CAR model, and we observe similar trends with our time-domain HiFi-GAN as well, which we include in the supplementary website linked in Section 1.

Figure 2 contains the mel-spectrograms of the generated speech from our model for each of these combined articulatory features. Our model is able to generalize to the unseen articulatory features between the three sounds. Moreover, 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.

### 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 methodology as our vowel interpolation experiment in Section 6.1.

Figure 3 depicts the mel-spectrograms of synthesized interpolation samples from our time-domain articulatory synthesis model. Similarly to our vowel interpolation results, we observe that our model generalizes to the unseen samples between the three consonants and exhibits smooth generation. Specifically, these results indicate that our model can smoothly transition between nasal, approximant, and lateral approximant consonants, similarly to the human speech production process.

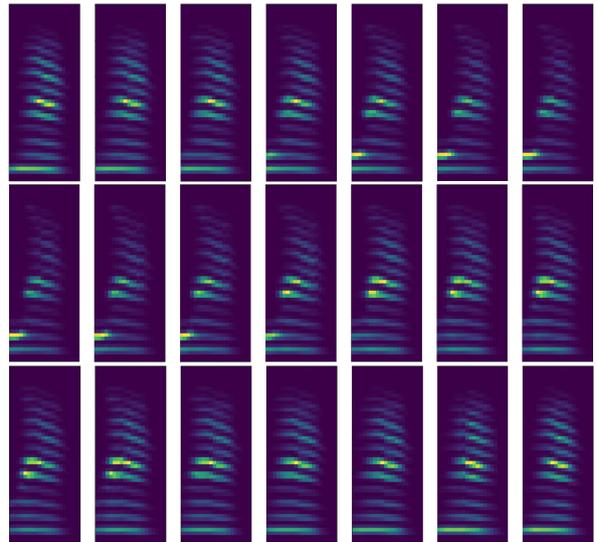


Figure 3: Alveolar consonant interpolation. The top row contains the synthesized samples between the "ra" and "na" sounds, the middle row "na" and "la", and the bottom row "la" and "ra".

To study how well our model generalizes across place of articulation, we fix the consonant type and interpolate between two places. Namely, we interpolate between the approximant consonants "ra" and "ja", using the same aforementioned methodology. Figure 4 depicts these results. As with our alveolar consonant interpolation results, we observe that our model generalizes to unseen samples and produces smooth transitions between synthesized interpolation samples here.

To quantify how the synthesized utterances

change across the interpolation, we create two plots studying changes in the magnitudes of different bands of the mel-spectrogram. Namely, our first graph plots the magnitude of each mel-spectrogram frequency vector across the seven utterances, going left to right in the interpolation. Our second plot does the same with time vectors, i.e., columns in the mel-spectrograms. We compute the magnitude of each vector using the L1 norm, which is just a sum here since mel-spectrogram values are non-negative. To improve readability in both plots, we omit vectors that on average change less than 0.3 in magnitude between adjacent interpolation samples.

As shown in the bottom row of Figure 4, the vector magnitude lines are generally monotonic and almost linear in many cases when going left to right in the interpolation. This supports our hypothesis that our model has learnt to transition smoothly between consonants when synthesizing articulatory features.

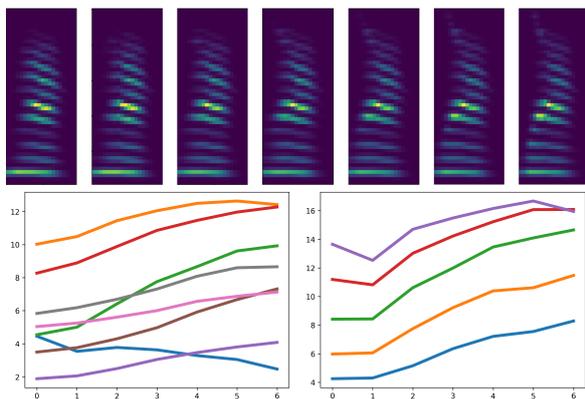


Figure 4: Approximate consonant interpolation. *Top row*: synthesized samples between the "ra" and "ja" sounds. *Bottom row left*: frequency vector magnitudes for each spectrum. *Bottom row right*: time vector magnitudes for each spectrum.

### 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 "tɛ" 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.

## 7 Synthesis Quality

### 7.1 Fidelity

Since MCD serves as an objective measure of synthesis quality (Black, 2019), we first measure synthesis fidelity using this metric. As detailed in Table 3, we observe that our time-domain articulatory synthesis approach achieves performance comparable to the frequency-domain baseline. Namely, our approach performs noticeably better than the baseline on the synthetic articulatory dataset and slightly worse on the EMA-to-speech task. Given these results, 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, which can be heard in the accompanying samples. We discuss this phoneme confusion in more detail below.

Model	MCD	
	Birk.-Pseudo	EMA-MGNU0
NSF-CAR	$3.36 \pm 0.28$	$5.44 \pm 0.67$
HiFi-GAN	<b><math>2.90 \pm 0.22</math></b>	$4.81 \pm 0.76$
Spec.-Int.	$5.15 \pm 0.48$	<b><math>4.75 \pm 0.81</math></b>

Table 3: MCD for each model on Birkholz and EMA data.

### 7.2 Automatic Speech Recognition

To evaluate the intelligibility of our synthesis approach, we conduct open-vocabulary transcription experiments for the EMA-to-speech task with our time-domain HiFi-GAN model described in Section 3.3. First, we perform an objective evaluation using deep automatic speech recognition (ASR) models. Specifically, we use DeepSpeech<sup>1</sup> (Hannun et al., 2014) as done by Gaddy and Klein (Gaddy and Klein, 2021) as well as the ESPnet Conformer ASR model trained on LibriSpeech<sup>2</sup> (Guo et al., 2021; Panayotov et al., 2015). We use these models to transcribe the synthesis outputs of our model on the entire MNGU0 evaluation set described in Section 4.1 and calculate the average word error rates (WERs) and character error rates (CERs). Since some utterances in the evaluation set contain proper nouns, we also compute ASR

<sup>1</sup><https://github.com/mozilla/DeepSpeech>

<sup>2</sup><https://zenodo.org/record/4604066#.YeNA0i2z2CM>

metrics on all of the evaluation set utterances composed entirely of common nouns, which form a 32-utterance subset.

Table 4 summarizes our ASR results. On the common-noun subset, our model achieves a character error rate of 10.7% with the ESPnet ASR model, indicating that our model is able to synthesize intelligible speech. The consistent differences between the WER and CER values as well as the entire set and common-noun subset performances suggests that these ASR metrics may be underestimating intelligibility, as also observed by Gaddy and Klein (Gaddy and Klein, 2021). Thus, we also evaluate the intelligibility of our model through human evaluations, as discussed in Section 7.3 below.

ASR Model	WER		CER	
	All	Com.	All	Com.
ESPnet	32.9	19.2	17.9	10.7
DeepSpeech	41.3	32.9	20.2	15.5

Table 4: ASR. entire evaluation set (All) and common noun subset (Com.).

### 7.3 Human Evaluation

To further understand the intelligibility of our time-domain articulatory synthesis approach, we also perform open-vocabulary transcription tests with human listeners, evaluating our same time-domain HiFi-GAN model (Section 3.3) used in our Section 7.2 ASR experiments above. Namely, we randomly select ten utterances from our EMA corpus evaluation set, choosing among the 32 sentences without proper nouns. Based on the transcriptions from six English-speaking listeners, our model achieves an average WER of 7.14%, indicating that our model is able to produce intelligible speech. To our knowledge, this value is noticeably lower than prior results, which are around 30.1% (Taguchi and Kaburagi, 2018). This suggests that our time-domain articulatory synthesis methodology is a suitable approach for efficiently performing speech synthesis while achieving high intelligibility.

## 8 Phoneme Confusion

To further study the phonological errors made by our model, we analyze the phonemes that our EMA-to-speech model confused during synthesis. Namely, we study phoneme confusability for our time-domain HiFi-GAN model (Section 3.3) through the transcriptions, both from the ASR ones

described in Section 7.2 and the human ones described in Section 7.3. For each transcribed utterance, we convert the graphemes to a phoneme sequence using Phonemizer<sup>3</sup> (Bernard and Titeux, 2021) and their eSpeak NG backend,<sup>4</sup> and repeat this grapheme-to-phoneme conversion with the ground truth texts. We identify the phoneme confusion pairs using sclite,<sup>5</sup> which aligns each predicted sequence with the respective ground truth and then records the substitution errors.

For our human evaluation analysis, we use all of the transcripts from the six listeners, i.e., 60 utterances. Figure 5 depicts the resulting phoneme confusion pairs. We plot these confusion pairs on an International Phonetic Alphabet (IPA) chart that extends the one from Gaddy and Klein to more phonemes (Gaddy and Klein, 2021), indicating pairs with a higher frequency of substitution errors using darker lines. We also populate this IPA chart with our confusion pairs from the ASR transcriptions in Figure 6, for which we use the texts transcribed by the ESPnet model for the entire MNGU0 evaluation set, as discussed in Section 7.2. We omit the phoneme pairs that are only confused once in Figure 6 in order to improve readability.

From these two IPA charts, we observe that the most of the word substitution errors are due to plosive or vowel confusions. Since the primary vowel confusions in Figure 5 differ from those in Figure 6, we hypothesize that vowel confusability for human evaluators mainly resulted from the substitution of vowels to form logical, grammatically correct words and phrases. The automatic transcribers may not have as much of such bias and we observe that the primary confused vowel pairs are relatively close to each other with our ASR-based results, reinforcing this hypothesis. One potential reason for the plosive substitutions is that plosives generally have a shorter duration than other consonant types like fricatives (Alwan et al., 2011) and thus may be more readily confusable. Among the plosives, "p", "b", "t", and "d" may have been easier to confuse than "k" and "g" for the human evaluators because the latter two plosives have longer voice onset times, a pattern also observed by Birkholz (Birkholz, 2013b). From Figure 6, we also observe that multiple voiced-unvoiced pairs are confused. We hypothesize that this is because the only voic-

<sup>3</sup><https://github.com/bootphon/phonemizer>

<sup>4</sup><https://github.com/espeak-ng/espeak-ng>

<sup>5</sup><https://github.com/usnistgov/SCTK>

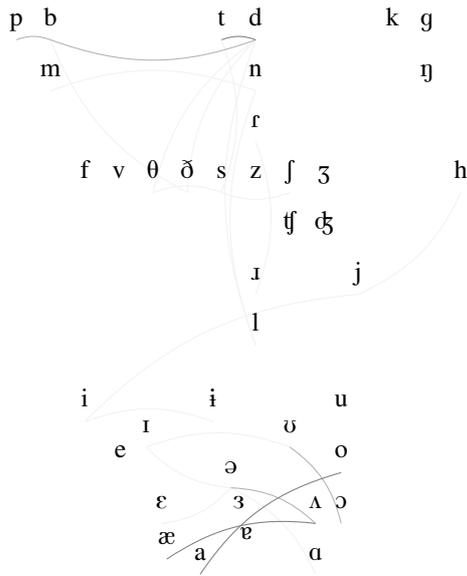


Figure 5: Phoneme confusability based on human transcriptions. Phoneme pairs that are confused more frequently have darker lines.

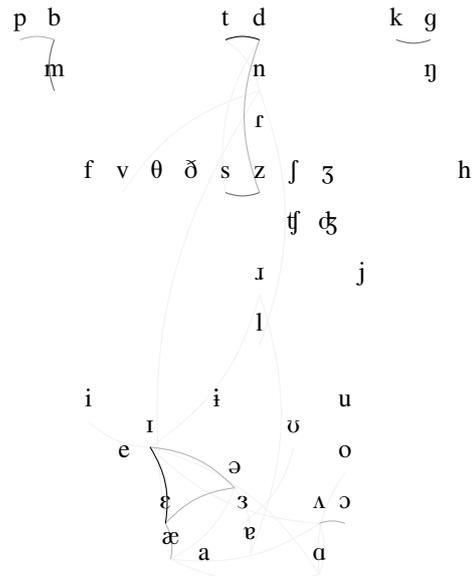


Figure 6: Phoneme confusability based on ASR transcriptions. Phoneme pairs that are confused more frequently have darker lines.

ing information that our EMA-to-speech model receives as input is the estimated F0 sequence, as described in Section 4.1.

## 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 that our model is also highly intelligible, achieving a transcription word error rate (WER) of 7.14% for the EMA-to-speech task. Moving forward, we plan to test our methodology on other modalities like electromyography (EMG) (Gaddy and Klein, 2021) and real-time magnetic resonance imaging (RT-MRI) (Lim et al., 2021). We also plan to extend our approach to multi-speaker and multilingual settings (Richmond et al., 2011; Lim et al., 2021; Wu et al., 2021b).

## References

Abeer Alwan, Jintao Jiang, and Willa Chen. 2011. Perception of place of articulation for plosives and fricatives in noise. *Speech communication*, 53(2):195–209.

Miguel Angrick, Christian Herff, Emily Mugler, Matthew C Tate, Marc W Slutzky, Dean J Krusienski, and Tanja Schultz. 2019. Speech synthesis from ecog using densely connected 3d convolutional neural networks. *Journal of neural engineering*, 16(3):036019.

Gopala K Anumanchipalli, Josh Chartier, and Edward F Chang. 2019. Speech synthesis from neural decoding of spoken sentences. *Nature*, 568(7753):493–498.

Sandesh Aryal and Ricardo Gutierrez-Osuna. 2016. Data driven articulatory synthesis with deep neural networks. *Computer Speech & Language*, 36:260–273.

Rohan Badlani, Adrian Łancucki, Kevin J Shih, Rafael Valle, Wei Ping, and Bryan Catanzaro. 2021. One tts alignment to rule them all. *arXiv preprint arXiv:2108.10447*.

Mathieu Bernard and Hadrien Titeux. 2021. **Phonemizer: Text to phones transcription for multiple languages in python**. *Journal of Open Source Software*, 6(68):3958.

Peter Birkholz. 2013a. Modeling consonant-vowel coarticulation for articulatory speech synthesis. *PloS one*, 8(4):e60603.

Peter Birkholz. 2013b. **Modeling consonant-vowel coarticulation for articulatory speech synthesis**. *PloS one*, 8:e60603.

Alan W Black. 2019. CMU wilderness multilingual speech dataset. In *ICASSP*, pages 5971–5975. IEEE.

619	Florent Bocquelet, Thomas Hueber, Laurent Girin, Pierre Badin, and Blaise Yvert. 2014. Robust articulatory speech synthesis using deep neural networks for bci applications. In <i>15th Annual Conference of the International Speech Communication Association (Interspeech 2014)</i> .	Takamichi, and Shinji Watanabe. 2021. Espnet2-tts: Extending the edge of tts research. <i>arXiv preprint arXiv:2110.07840</i> .	675 676 677
625	Yu-Wen Chen, Kuo-Hsuan Hung, Shang-Yi Chuang, Jonathan Sherman, Wen-Chin Huang, Xugang Lu, and Yu Tsao. 2021. Ema2s: An end-to-end multimodal articulatory-to-speech system. In <i>2021 IEEE International Symposium on Circuits and Systems (ISCAS)</i> , pages 1–5. IEEE.	Hirofumi Inaguma, Shun Kiyono, Kevin Duh, Shigeki Karita, Nelson Yalta, Tomoki Hayashi, and Shinji Watanabe. 2020. <a href="#">ESPnet-ST: All-in-one speech translation toolkit</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 302–311, Online. Association for Computational Linguistics.	678 679 680 681 682 683 684
631	Tam'as G'abor Csap'o, Csaba Zaink'o, L. Viktor T'oth, Gábor Gosztolya, and Alexandra Mark'o. 2020. Ultrasound-based articulatory-to-acoustic mapping with waveglow speech synthesis. In <i>Interspeech</i> .	Khalil Iskarous, Louis Goldstein, Douglas H Whalen, Mark Tiede, and Philip Rubin. 2003. Casy: The haskins configurable articulatory synthesizer. In <i>International Congress of Phonetic Sciences, Barcelona, Spain</i> , pages 185–188.	685 686 687 688 689
635	Isaac Elias, Heiga Zen, Jonathan Shen, Yu Zhang, Jia Ye, R. J. Skerry-Ryan, and Yonghui Wu. 2021. Parallel tacotron 2: A non-autoregressive neural tts model with differentiable duration modeling. <i>ArXiv</i> , abs/2103.14574.	Ye Jia, Michelle Tadmor Ramanovich, Tal Remez, and Roi Pomerantz. 2021. Translatotron 2: Robust direct speech-to-speech translation. <i>arXiv preprint arXiv:2107.08661</i> .	690 691 692 693
640	Gunnar Fant. 1991. What can basic research contribute to speech synthesis? <i>Journal of Phonetics</i> , 19(1):75–90.	Ye Jia, Ron Weiss, Fadi Biadisy, Wolfgang Macherey, Melvin Johnson, Zhifeng Chen, and Yonghui Wu. 2019. <a href="#">Direct speech-to-speech translation with a sequence-to-sequence model</a> . In <i>Interspeech</i> , pages 1123–1127.	694 695 696 697 698
643	Gunnar Fant. 1995. The lf-model revisited. transformations and frequency domain analysis. <i>Speech Trans. Lab. Q. Rep., Royal Inst. of Tech. Stockholm</i> , 2(3):40.	A Karmel, Anushka Sharma, Muktak pandya, and Diksha Garg. 2019. <a href="#">Iot based assistive device for deaf, dumb and blind people</a> . <i>Procedia Computer Science</i> , 165:259–269. 2nd International Conference on Recent Trends in Advanced Computing ICRTAC -DISRUP - TIV INNOVATION , 2019 November 11-12, 2019.	699 700 701 702 703 704 705
646	David Gaddy and Dan Klein. 2021. <a href="#">An improved model for voicing silent speech</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 175–181, Online. Association for Computational Linguistics.	Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. <a href="#">Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech</a> . In <i>Proceedings of the 38th International Conference on Machine Learning</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 5530–5540. PMLR.	706 707 708 709 710 711
653	Marc-Antoine Georges, Pierre Badin, Julien Diard, Laurent Girin, Jean-Luc Schwartz, and Thomas Hueber. 2020. Towards an articulatory-driven neural vocoder for speech synthesis. In <i>International Seminar on Speech Production</i> .	Jong Wook Kim, Justin Salamon, Peter Li, and Juan Pablo Bello. 2018. Crepe: A convolutional representation for pitch estimation. In <i>2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 161–165. IEEE.	712 713 714 715 716
658	Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jiatong Shi, et al. 2021. Recent developments on espnet toolkit boosted by conformer. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5874–5878. IEEE.	Diederik P. Kingma and Jimmy Ba. 2015. <a href="#">Adam: A method for stochastic optimization</a> . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	717 718 719 720 721
666	Awni Y. Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Gregory Frederick Diamos, Erich Elsen, Ryan J. Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, and A. Ng. 2014. Deep speech: Scaling up end-to-end speech recognition. <i>ArXiv</i> , abs/1412.5567.	Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. 2020. <a href="#">Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 17022–17033. Curran Associates, Inc.	722 723 724 725 726 727
672	Tomoki Hayashi, Ryuichi Yamamoto, Takenori Yoshimura, Peter Wu, Jiatong Shi, Takaaki Saeki, Yooncheol Ju, Yusuke Yasuda, Shinnosuke	Kevin Lenzo and Alan Black. 2000. Diphone collection and synthesis. <i>ICSLP</i> .	728 729

730	Yongwan Lim, Asterios Toutios, Yannick Bliesener,	Simon Stone, Philipp Schmidt, and Peter Birkholz. 2020.	782
731	Ye Tian, Sajan Lingala, Colin Vaz, Tanner Sorensen,	Prediction of voicing and the f0 contour from elec-	783
732	Miran Oh, Sarah Harper, Weiyi Chen, Yoonjeong	tromagnetic articulography data for articulation-to-	784
733	Lee, Johannes Töger, Mairym Llorens Monteserin,	speech synthesis. In <i>ICASSP 2020-2020 IEEE Inter-</i>	785
734	Caitlin Smith, Bianca Godinez, Louis Goldstein,	<i>national Conference on Acoustics, Speech and Signal</i>	786
735	Dani Byrd, Krishna Nayak, and Shrikanth Narayanan.	<i>Processing (ICASSP)</i> , pages 7329–7333. IEEE.	787
736	2021. <a href="#">A multispeaker dataset of raw and recon-</a>		
737	<a href="#">structed speech production real-time mri video and</a>	Fumiaki Taguchi and Tokihiko Kaburagi. 2018.	788
738	<a href="#">3d volumetric images.</a> <i>Scientific Data</i> , 8.	Articulatory-to-speech conversion using bi-	789
		directional long short-term memory. In <i>Interspeech</i> ,	790
739	Zheng-Chen Liu, Zhen-Hua Ling, and Li-Rong Dai.	pages 2499–2503.	791
740	2018. Articulatory-to-acoustic conversion using		
741	blstm-rnns with augmented input representation.	Andros Tjandra, Sakriani Sakti, and Satoshi Nakamura.	792
742	<i>Speech Communication</i> , 99:161–172.	2019. <a href="#">Speech-to-speech translation between untran-</a>	793
		<a href="#">scribed unknown languages.</a> In <i>2019 IEEE Auto-</i>	794
743	Max Morrison, Rithesh Kumar, Kundan Kumar, Prem	<i>matic Speech Recognition and Understanding Work-</i>	795
744	Seetharaman, Aaron Courville, and Yoshua Bengio.	<i>shop (ASRU)</i> , pages 593–600.	796
745	2022. Chunked autoregressive gan for conditional		
746	waveform synthesis. In <i>Submitted to ICLR 2022</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	797
		Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz	798
747	Tomáš Nekvinda and Ondřej Dušek. 2020. <a href="#">One Model,</a>	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	799
748	<a href="#">Many Languages: Meta-Learning for Multilingual</a>	<a href="#">you need.</a> In <i>Advances in Neural Information Pro-</i>	800
749	<a href="#">Text-to-Speech.</a> In <i>Interspeech</i> , pages 2972–2976.	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	801
		Xin Wang, Shinji Takaki, and Junichi Yamagishi. 2019.	802
750	Vassil Panayotov, Guoguo Chen, Daniel Povey, and	<a href="#">Neural source-filter-based waveform model for sta-</a>	803
751	Sanjeev Khudanpur. 2015. Librispeech: an asr corpus	<a href="#">tistical parametric speech synthesis.</a> In <i>ICASSP 2019</i>	804
752	based on public domain audio books. In <i>2015</i>	<i>- 2019 IEEE International Conference on Acoustics,</i>	805
753	<i>IEEE international conference on acoustics, speech</i>	<i>Speech and Signal Processing (ICASSP)</i> , pages 5916–	806
754	<i>and signal processing (ICASSP)</i> , pages 5206–5210.	5920.	807
755	IEEE.		
		Yuxuan Wang, R. J. Skerry-Ryan, Daisy Stanton,	808
756	Adam Polyak, Yossi Adi, Jade Copet, Eugene	Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, Zongheng	809
757	Kharitonov, Kushal Lakhotia, Wei-Ning Hsu, Ab-	Yang, Ying Xiao, Z. Chen, Samy Bengio, Quoc V. Le,	810
758	delrahman Mohamed, and Emmanuel Dupoux. 2021.	Yannis Agiomyrgiannakis, Robert A. J. Clark, and	811
759	Speech Resynthesis from Discrete Disentangled Self-	Rif A. Saurous. 2017. Tacotron: Towards end-to-end	812
760	Supervised Representations. In <i>Interspeech</i> .	speech synthesis. In <i>Interspeech</i> .	813
		Peter Wu, Paul Pu Liang, Jiatong Shi, Ruslan Salakhut-	814
761	Ryan Prenger, Rafael Valle, and Bryan Catanzaro. 2019.	dinov, Shinji Watanabe, and Louis-Philippe Morency.	815
762	<a href="#">Waveglow: A flow-based generative network for</a>	2021a. Understanding the tradeoffs in client-side	816
763	<a href="#">speech synthesis.</a> In <i>ICASSP 2019 - 2019 IEEE Inter-</i>	privacy for downstream speech tasks. In <i>APSIPA</i>	817
764	<i>national Conference on Acoustics, Speech and Signal</i>	<i>ASC</i> .	818
765	<i>Processing (ICASSP)</i> , pages 3617–3621.		
		Peter Wu, Jiatong Shi, Yifan Zhong, Shinji Watanabe,	819
766	Korin Richmond, Phil Hoole, and Simon King. 2011.	and Alan W Black. 2021b. Cross-lingual transfer for	820
767	<a href="#">Announcing the electromagnetic articulography (day</a>	speech processing using acoustic language similarity.	821
768	<a href="#">1) subset of the mngu0 articulatory corpus.</a> In <i>Inter-</i>	In <i>ASRU</i> .	822
769	<i>speech</i> , pages 1505–1508.		
		Chengzhu Yu, Heng Lu, Na Hu, Meng Yu, Chao Weng,	823
770	Philip Rubin, Thomas Baer, and Paul Mermelstein.	Kun Xu, Peng Liu, Deyi Tuo, Shiyin Kang, Guangzhi	824
771	1981. An articulatory synthesizer for perceptual re-	Lei, et al. 2019. Durian: Duration informed attention	825
772	search. <i>The Journal of the Acoustical Society of</i>	network for multimodal synthesis. <i>arXiv preprint</i>	826
773	<i>America</i> , 70(2):321–328.	<i>arXiv:1909.01700</i> .	827
		Yu Zhang, Ron J Weiss, Heiga Zen, Yonghui Wu,	828
774	Celia Scully. 1990. Articulatory synthesis. In <i>Speech</i>	Zhifeng Chen, RJ Skerry-Ryan, Ye Jia, Andrew	829
775	<i>production and speech modelling</i> , pages 151–186.	Rosenberg, and Bhuvana Ramabhadran. 2019. Learn-	830
776	Springer.	ing to speak fluently in a foreign language: Multi-	831
		lingual speech synthesis and cross-language voice	832
777	Berrak Sisman, Junichi Yamagishi, Simon King, and	cloning. <i>arXiv preprint arXiv:1907.04448</i> .	833
778	Haizhou Li. 2020. An overview of voice conversion		
779	and its challenges: From statistical modeling to deep	<b>A Appendix</b>	834
780	learning. <i>IEEE/ACM Transactions on Audio, Speech,</i>		
781	<i>and Language Processing</i> .		

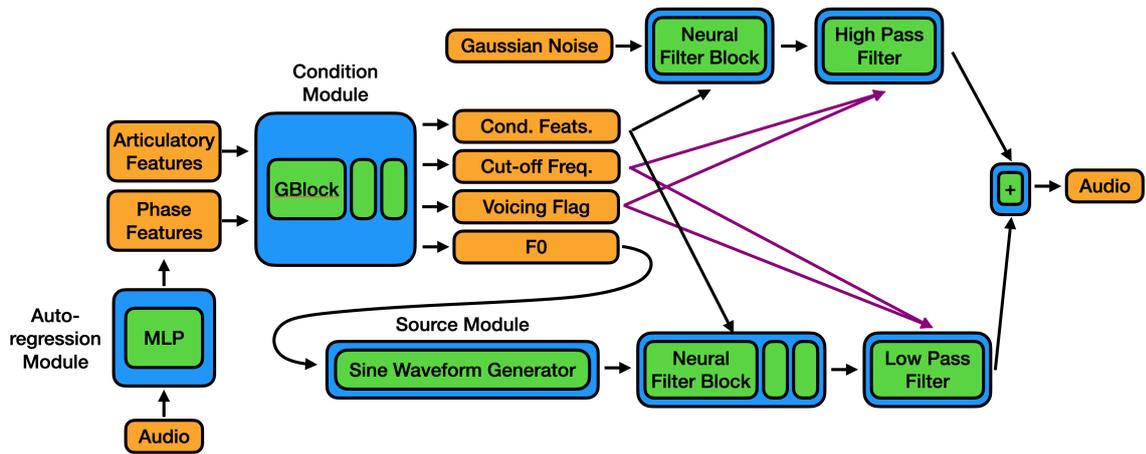


Figure 7: Model architecture of our NSF-CAR generator.