

Multimodal Pretraining for Vocal Tract Modeling

Anonymous CVPR submission

Paper ID 14145

Abstract

Accurate modeling of the vocal tract is necessary for naturalistic facial animation, avatar rendering for virtual environments, and pronunciation tutoring. However, vocal tract modeling is challenging because internal articulators like the tongue and velum are occluded from external motion capture technologies. Real-time magnetic resonance imaging (RT-MRI) offer a direct way to measure precise movements of internal articulators during naturalistic speaking, offering a possible solution to modeling. However, segmented and annotated datasets of MRI are limited in size due to time-consuming and computationally expensive labeling methods. We first present a labeling strategy for the Speech MRI Open Dataset comprising over 20 hours of audio-aligned RT-MRI video using a vision-only segmentation approach. We then apply a multimodal pre-training algorithm to: (1) improve segmentation of vocal articulators, (2) synthesize intelligible speech from inferred segments, and (3) animate 2D and 3D facial avatars that capture the complex articulator patterns underlying naturalistic speech production. We also extend these animation techniques to a high performance streamable facial avatar driven directly from speech, achieving 107ms average latency. Together, we set a new benchmark for vocal tract modeling in MRI image segmentation, intelligible MRI-to-speech synthesis, and real-time speech-to-avatar rendering.

1. Introduction

Vocal tract modeling is an essential technology in many applications including facial animation, naturalistic speaking avatars, and second language pronunciation learning [21, 29, 34, 50]. Modeling is also necessary in healthcare applications such as Brain-Computer Interfaces for communication [5] and diagnosing and treating speech disfluencies [32, 40].

Methods of external motion capture cannot record precise and accurate vocal tract movements for occluded articulators. Hence, the inner mouth is often poorly lit or neglected in multimedia approaches to motion capture-based

facial animation [33].

Popular approaches to solving the issue of inner mouth occlusion include electromagnetic articulography (EMA) and electromyography (EMG) as models for the vocal tract. However, these methods only contain a small subset of articulatory features.

A more comprehensive approach uses Real-Time Magnetic Resonance Imaging (RT-MRI) of the vocal tract. This technology offers audio-aligned videos of internal and external articulators that are not measurable by other articulatory representations. When tested against downstream speech-related tasks, RT-MRI has been shown to more reliably and completely model the vocal tract in comparison to EMA [59]. However, current state-of-the-art labeling methods for extracting interpretable features from these videos are time-consuming, computationally expensive, and prone to errors [8]. Therefore, only a small amount of vocal tract RT-MRI data is labeled [35].

In this paper, we propose a comprehensive application of pretraining with both video and audio modalities for modeling the vocal tract. We first present a high-performance vision-based and multimodal RT-MRI feature extraction approach. Using these results, we label the Speech MRI Open Dataset [30] containing over 20 hours of vocal tract RT-MRI data for 75 speakers diverse in age, gender, and accent. To our knowledge, this dataset increases the amount of labeled public RT-MRI data of the vocal tract by over a factor of 9.

Using this newly labeled dataset, we deploy an MRI-EMA pretraining method to further evaluate our feature extraction models using MRI-based deep speech synthesis. We achieve significant improvements in intelligibility compared to synthesis using the ground truth MRI tracks.

Another downstream application of our multimodal pre-training is 2D and 3D facial avatar visualization. We propose a deep articulatory inversion technique for speech-to-MRI prediction with a direct mapping to the 3D avatar, enabling us to visualize complex naturalistic speech production. We employ this method in both offline and real-time scenarios, achieving an average streaming latency of 107ms/batch of 100ms audio data. This result is a 20% improvement over the previous baseline despite using a 3D

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078

079 model with more than $6 \times$ the vertices.
080

081 To briefly summarize, our contributions in this work
082 are: (1) a labeled version of the 75-speaker Speech MRI
083 Open Dataset, (2) the application of a multimodal pretrain-
084 ing method for RT-MRI analysis, and (3) low latency facial
085 avatar visualizations of MRI-based vocal tracts in offline
086 and real-time contexts. These architectures are outlined in
Figure 1.

087 2. Related Work

088 2.1. MRI Feature Extraction

089 Most published work using vocal tract RT-MRI feature ex-
090 tractions either do not use high-dimensional interpretable
091 vocal tract representations, or only use previously extracted
092 features for downstream tasks without generalizing well to
093 unseen speakers [22, 32, 59, 66]. The existing algorithm for
094 speaker segmentation traces airway contours using hand-
095 drawn reference boundaries. It extracts 170 points of MRI
096 coordinates per frame, taking up to 20 minutes to converge
097 for a single frame [35]. From this point forward, we re-
098 fer to the outputs of this algorithm as the “ground truth.”
099 More recent work in [2] uses a deep attention-gated U-Net
100 for multi-speaker deep RT-MRI feature extraction, which is
101 more efficient and generalizes well to unseen speakers.

102 2.2. Deep MRI-Based Speech Synthesis

103 Our concurrent submission attached in supplemental mate-
104 rials trains single-speaker MRI to speech. We detail our ap-
105 proach in Section 5.2. In this work, we explore new MRI-to-
106 speech methodologies using the newly labelled 75-speaker
107 dataset.

108 2.3. Speech-Driven Avatar Animation

109 Within the domain of automated facial animation, there
110 have been many different approaches to driving an avatar
111 from speech. Linguistic methods aim to map phonemes
112 to visemes on the avatar [16, 44, 52, 63]. However, these
113 systems require manually defined complex rules without a
114 streamlined approach for inner mouth animation. Other ap-
115 proaches to speech-driven animation include deep learning
116 models that are trained on audio-mesh paired data to per-
117 form mesh deformations [6, 15, 62]. However, these meth-
118 ods typically do not accurately model the tongue or suffer
119 from oversmoothing when directly regressing to facial mesh
120 movements [62].

121 Inner mouth animation from medical imaging represen-
122 tations of the vocal tract has been explored in [11], which
123 provides real-time (21 Hz) 3D tongue animation using a
124 streaming ultrasound snake contour extraction algorithm.
125 Due to the tongue tip not being captured well in ultrasound,
126 this tongue model is incomplete. Similarly, in [49], a kine-
127 matic tongue model extracted from a single MRI volumetric

128 scan was directly animated using offline EMA data. These
129 works demonstrate the benefits of being grounded within
130 physiology but remain inaccessible to users without access
131 to equipment for directly capturing the vocal tract.

132 Advances in deep articulatory inversion models have
133 demonstrated the ability to approximate the physiology-
134 grounded representations of the vocal tract from solely
135 speech inputs [7, 9, 10, 20, 24–26, 28, 31, 37, 41, 43, 45,
136 47, 48, 51, 54, 55, 57, 58, 60, 61, 64, 65]. These approaches
137 avoid potentially invasive recording equipment while re-
138 taining the valuable vocal tract information from speech
139 production.

140 Our concurrent submission included in supplementary
141 materials [3] builds on these works. We use WavLM [12]
142 as feature extractors for input speech with multitask learn-
143 ing to predict tract variables, phonemes, pitch, and EMA
144 simultaneously. In this work, we use a similar multi-task
145 learning approach but with a pretraining RT-MRI and EMA
146 inversion model on 75 RT-MRI speakers and 8 HPRC EMA
147 speakers as described in Section 4.2 [53].

148 More recently, Medina *et al.* [33] combines the two
149 methods of deep articulatory inversion models and vo-
150 cal tract image representations. The result is a high-
151 performance offline solution to facial animation with an em-
152 phasis on the inner mouth using EMA as an intermediate
153 physiological representation. However, this approach has
154 a low-dimensional tongue rig and loses physiological in-
155 formation when optimizing for the Metahuman FACS rig
156 [17, 19].

157 Inspired by this speech-driven system, we developed a
158 streaming system for speech-to-avatar synthesis as a con-
159 current paper submission (included in supplementary ma-
160 terials) [1]. This work builds on the approach from our
161 previous submission in the following ways: (1) introducing
162 RT-MRI as a replacement intermediary feature for EMA
163 to increase inner mouth mesh resolution, (2) grounding the
164 3D facial model within human physiology according to RT-
165 MRI of the vocal tract, and (3) optimizing streaming us-
166 ing a new system in Unreal Engine as opposed to Autodesk
167 Maya.

168 3. Datasets

169 3.1. USC-TIMIT Dataset

170 We use the labeled 8-speaker RT-MRI USC-TIMIT dataset
171 of the vocal tract for training and as the ground truth feature
172 extractions for all baseline and experimental approaches de-
173 scribed in [35]. Given sentences on a projection screen, sub-
174 jects were instructed to read each out at a natural speaking
175 rate while laying supine in an MRI scanner. A four-channel
176 upper airway receiver coil array was used for receiving sig-
177 nals, which were processed to reproduce 84×84 midsag-
178 gital MR videos capturing lingual, labial, and jaw motion,

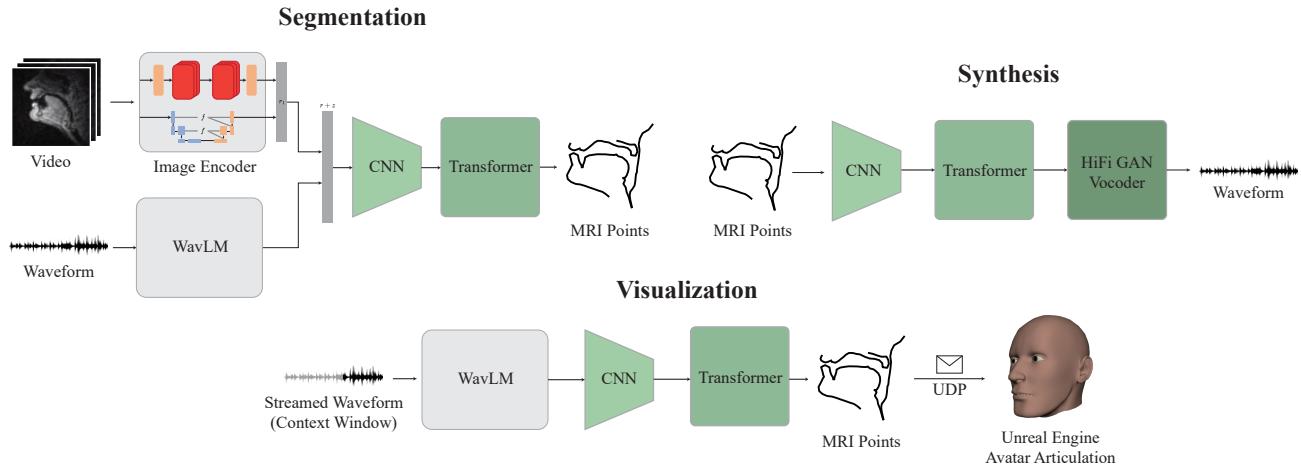


Figure 1. Summary of our approach showing three architectures for MRI analysis in segmentation, synthesis, and 3D visualization.

179 and velum, pharynx, and larynx articulations. These videos
 180 are collected at 83.33 Hz, and the algorithm described in
 181 Section 2.1 is used to extract the 170 representative points
 182 of the vocal tract. Of these 170 points, we take the subset
 183 of 95 points (190×2 coordinates) that has been deter-
 184 mined to be most vital for speech tasks in [59].

185 Paired with these trajectories is the 16kHz speech data
 186 (resampled from original 20kHz) corresponding to the read
 187 sentence during any RT-MRI scan. We further enhanced
 188 this audio using Adobe Podcast to reduce reverberation, as
 189 done in [59].

190 3.2. Speech MRI Open Dataset

191 The Speech MRI Open Dataset [30] is a multispeaker
 192 dataset from USC that provides synchronized speech of 75
 193 diverse speakers with raw multi-coil RT-MRI videos of the
 194 vocal tract during articulation. Such a large, rich dataset
 195 can help solve many open problems in fields related to
 196 phonetics, spoken language, and vocal articulation. However,
 197 unlike the USC-TIMIT dataset, the data does not include
 198 labeled MRI feature points tracked over time, except for 6
 199 speakers. This labeled data is not yet available to the pub-
 200 lic.

201 4. Multimodal Feature Extraction

202 4.1. Image-Based Feature Extraction

203 We follow [2], a U-Net style model [42] with an atten-
 204 tion gating mechanism [36], for a baseline in MRI fea-
 205 ture extraction using the frames alone as the input modal-
 206 ity. We learn a spatial weighting map for each predicted
 207 point trained using Kullback-Leibler (KL) divergence loss
 208 between the weighting map and a 2D Gaussian heatmap.
 209

This model is trained on 66 minutes of ground truth RT-

210 MRI data from 7 of the 8 USC-TIMIT speakers. Given an
 211 MRI image, the U-Net predicts 95 84×84 spatial weight-
 212 ing maps for each of the 95 MRI points, and converts to points
 213 using a weighted average of the 25 highest pixels.

214 Our concurrent submission's [2] U-Net (in supple-
 215 mentary materials) is being used in our work as a baseline label-
 216 ing model that informs future models.

217 An additional challenge with the ground truth data is the
 218 presence of jitter, or random high-frequency perturbations
 219 across consecutive frames, leading to noisy point tracking.
 220 While the trained U-Net outputs smoother point tracks, we
 221 additionally apply a temporal Gaussian low-pass filter inde-
 222 pendently for each point. This results in far smoother tracks,
 223 while maintaining accurate point detection per frame.

224 Using the U-Net model as a pretrained convolutional input,
 225 we further explore joint point tracking using a convolu-
 226 tional LSTM as in [22] (CLSTM) and a Transformer. The
 227 CLSTM, previously used in MRI video segmentation [66],
 228 applies a 2-layer LSTM to the predicted U-Net outputs,
 229 trained on speech from the same 7 USC-TIMIT speakers.
 230 The Transformer similarly uses the U-Net points from each
 231 timestep, with an additional positional encoding. Traditionally,
 232 multi-frame point tracking is done using optical flow
 233 [14] or by extension, Kalman filtering [13]. Recent joint
 234 point tracking results have found that explicitly appending
 235 optical flow to the Transformer input has resulted in bet-
 236 ter resulting tracks [46]. Thus, we further input predicted
 237 single frame optical flow, averaged over articulators using
 238 the Lucas-Kanade assumption that points in close prox-
 239 imity have similar flow in subsequent timesteps. Both the
 240 CRNN and the Transformer methods did not achieve equal
 241 or better performance than smoothed U-Net tracks on MRI
 242 videos of unseen speakers, reinforcing the fact that artic-
 243ulatory MRI tracking is fundamentally different than other

244 traditional video tracking problems. To address this, we ex-
245 plore additional modalities in the following sections.

246 4.2. Speech-Based Feature Extraction

247 For a secondary unimodal baseline in MRI feature extrac-
248 tion, we use the speech audio waveforms corresponding to
249 the USC-TIMIT MRI trajectories as the input modality. Us-
250 ing the 10th layer of WavLM, we derive speech representa-
251 tions from the audio as the input to a Transformer prepended
252 with three residual convolutional blocks.

253 Additionally, the Transformer model is trained on the
254 speech data from 7 of the 8 USC-TIMIT speakers for multi-
255 task learning and outputs MRI trajectories and pitch from
256 the speech representations simultaneously.

257 4.3. Multimodal Feature Extraction

258 We experiment with multiple multimodal models for fea-
259 ture extraction, using representations from video frames and
260 from speech waveforms. We concatenate the two represen-
261 tations as input to a Transformer. Following 4.2, we train
262 each of the multimodal models on the same 7 of 8 USC-
263 TIMIT speakers with weighted L1 loss on outputted MRI
264 and pitch.

265 4.4. Labeling Speech MRI Open Dataset

266 We deploy the previously described U-Net model trained
267 on data from 7 USC-TIMIT speakers and 5 of the 6 la-
268 beled Speech MRI Open Dataset speakers, with temporal
269 Gaussian low-pass filtering, to fully label video and audio
270 aligned MRI point trajectories for the entire Speech MRI
271 Open Dataset.

272 In Figure 2, we highlight the efficacy and generalizing
273 qualities of the U-Net model on unseen speakers, allowing
274 us to expand the amount of labeled MRI to over 20 hours
275 across 83 total speakers. Qualitatively, the predicted seg-
276 mentations closely follow the ground truth to trace the de-
277 sired MRI segments, achieving a high quality labeling for
278 unseen speakers. Further results are discussed in Section 7.
279 As part of this paper, we also present this labeling for use in
280 future downstream speech tasks, increasing the amount of
281 labeled articulatory RT-MRI data available by over a factor
282 of 9.

283 5. Deep MRI-Based Speech Synthesis

284 Another noteworthy challenge using labeled MRI data is
285 speech synthesis from the MRI articulatory space. This is
286 important for synchronized between avatar and speech ren-
287 dering. Due to the higher resolution of MRI compared to
288 the six points of EMA data used in prior works, a more de-
289 tailed study of the vocal tract during speech production can
290 be conducted. We explore single-speaker deep speech syn-
291 thesis using the newly-labeled Speech MRI Open Dataset

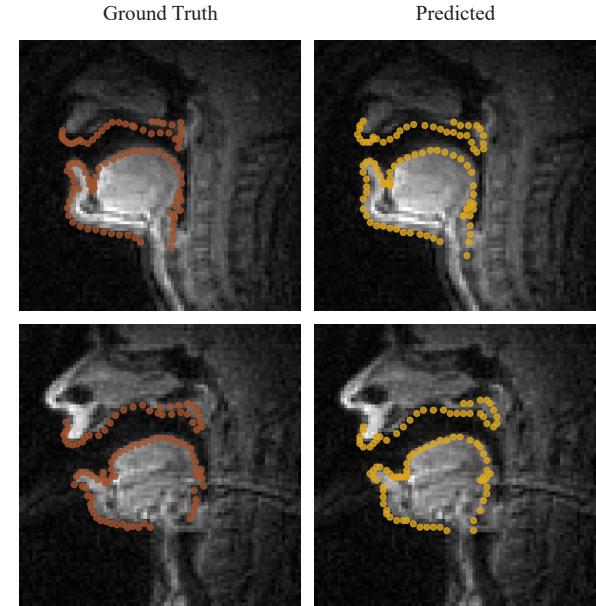


Figure 2. Two representative examples of predicted MRI points (right) compared to ground truth (left). Unseen examples spoken by unseen speakers from USC-TIMIT (top, Female) and Speech MRI Open Dataset (bottom, Male).

as an additional evaluation metric for potential future appli- 292
cations of this data. 293

294 5.1. Models

295 For our main speech synthesis architecture, we build on [4] 296 as illustrated in Figure 1. Given a set of 95 MRI point tra-
297 jectories, we follow [59] to first concatenate and flatten the
298 95 x - y pairs into a 190-length vector and center the data
299 around a point located on the hard palate with the lowest
300 standard deviation.

301 With this preprocessed MRI, we first use 3 1D convolu-
302 tional layers to encode the MRI points into an input for
303 a 6-layer Transformer [18, 56]. The output of the Trans-
304 former is 256-channel HuBERT vectors [23] finetuned on
305 VCTK with additional regularization loss. We then use a
306 HiFi-GAN vocoder [27] to directly synthesize speech from
307 these learned intermediate features. Full training details are
308 provided in the appendix.

309 5.2. Experiments

310 For evaluation, we pretrain the Transformer on a version of
311 the 75-speaker dataset labeled using a U-Net trained on 7
312 USC-TIMIT speakers. We combine the MRI dataset with
313 HPRC EMA dataset [53]. Our concurrent submission at-
314 tached in supplementary materials shows that pretraining
315 with both MRI and EMA modality helps the model gener-

316 alize to unseen examples. However, concurrent submission
317 only pretrained on single MRI speaker. With the 75-speaker
318 dataset, we also finetune this multi-modal pretrained trans-
319 former on 75-speaker MRI only. Given these two types
320 of pretrained weights, we finetune each of them on single
321 speaker data based on the chosen segmentation of a USC-
322 TIMIT speaker.

323 We choose the ground truth and the U-Net predicted
324 MRI trajectories of USC-TIMIT as two baseline metrics
325 to compare to the performance of the previously described
326 multimodal segmentation models. Additionally, we inde-
327 pendently finetune the model on two separate speakers: one
328 seen speaker and one unseen speaker (by the labeling U-
329 Net model). In this manner, we can assess the quality of
330 a given segmentation model on the important downstream
331 task of recovering intelligibility from the information-rich
332 MRI segmentations. Full model ablation is included in the
333 appendix.

334 6. Visualization

335 We further investigate the vocal tract during natural articu-
336 lation using offline and real-time 3D visualizations of the
337 face and mouth in combination with deep articulatory in-
338 version. We aim to provide an anatomically-coherent visual
339 representation of inferred MRI trajectories in speech pro-
340 duction. Building on [38], our animated 3D model provides
341 a faster, higher performance system for generating speech-
342 driven avatar movements.

343 6.1. Deep Articulatory Inversion

344 For the proposed vocal animation architecture, the first step
345 is to use an acoustic-to-articulatory inversion technique to
346 predict MRI trajectories that will feed a custom 3D fa-
347 cial model. We use the newly labeled data of 75 speakers
348 to learn to predict midsagittal x and y coordinates of the
349 tongue (20 points), hard palate, velum (15 points), lips (ℓ
350 points), lower incisor, and epiglottis.

351 We follow the same unimodal segmentation architec-
352 ture from Section 4.2 as a speech-to-MRI inversion model.
353 Combining the newly labeled 75-speaker MRI dataset with
354 the HPRC dataset, we first pretrain the model with two
355 heads, one outputting EMA, tract variables, phonemes, and
356 pitch, and the other head outputting MRI and pitch. We then
357 finetune the model on the 8-speaker TIMIT dataset with 1
358 speaker held out as test speaker.

359 6.2. Face and Vocal Tract Model

360 We construct a custom 3D face and mouth model in Au-
361 todesk Maya. We define point 88 of the 95-point MRI
362 trajectory subset on the hard palate to be the origin of the 3D
363 visualization space.

364 Relative to this point, we trace a tongue mesh according
365 to a frame from the USC-TIMIT Napa speaker. For each

366 tongue point in the MRI trajectories, we bind a joint to the
367 tongue model in Maya. We also bind a tongue centroid joint
368 for positional control of the 3D tongue and preservation of
369 the overall topology of the mesh during deformation.

370 We follow a similar procedure for the hard palate, lips,
371 and epiglottis to preserve movement trends and minimize
372 noise. Due to noise from the inversion model being un-
373 predictably amplified in high-resolution visualization of the
374 velum and lips, we infer a low-joint mapping of these fea-
375 tures from the original points from the MRI tracks. This
376 approach retains the general trend of the vocal tract move-
377 ments while disregarding the noise of individual MRI points
378 trajectories. Finally, we model the movement of the lower
379 incisor as a proxy for jaw movements through a hinge-based
380 approximation following [1].

381 The full face model is illustrated and labeled with im-
382 portant features in Figure 4. Additionally, in Figure 5, we
383 highlight a major benefit of using RT-MRI over EMA-based
384 representations of past works: we can support higher gran-
385 ularity in tongue movements through 20 joints instead of 3.

386 6.3. Offline Animation

387 For animating the vocal tract when given a full utterance, we
388 can offline the approach reported by [38] within Autodesk
389 Maya as a qualitative baseline. The deep articulatory inver-
390 sion model is deployed and provided the entire waveform as
391 input. The MRI trajectory outputs from this model are then
392 post-processed to correspond one-to-one with the rig joints
393 of the facial model in Maya. For each timestep of the MRI
394 trajectories, we transform each joint according to the MRI
395 track and either set a keyframe or refresh the Viewport 2.0
396 to simulate animation. In the keyframe case, we can sim-
397 ply then play the animation from the start frame to the end
398 frame at 83.33 Hz after all timesteps are keyed to achieve
399 offline vocal animation.

400 However, since scripting in Maya is a blocking process,
401 this method effectively freezes the program until animation
402 is complete. A user is unable to interact with the 3D model
403 and visually analyze speech production from multiple an-
404 gles. Additionally, using a keyframe-based approach en-
405 tails a delay before any animation is played while using
406 a viewport-based approach has high refresh latency as re-
407 ported by [38].

408 Alternatively, the real-time steps outlined in the follow-
409 ing sections can be adapted to resolve the aforementioned
410 issues. Deploying a network-based approach in an opti-
411 mized game engine alleviates both problems of interactivity
412 and latency.

413 6.4. Real-time Speech-to-MRI Processing

414 We build on and modify the general strategy proposed in [1]
415 to accommodate streamed audio for vocal animation. The
416 specific streaming system is briefly outlined below.

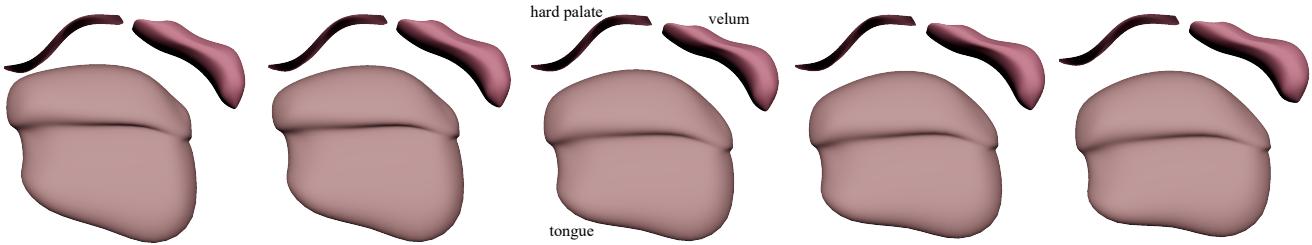


Figure 3. End of “buns”—Five frames of the tongue, hard palate, and velum generated for audio length of 60ms, corresponding to the fricative [s].

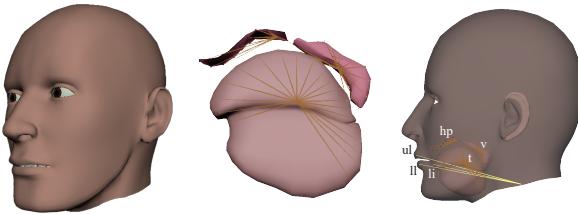


Figure 4. Three-quarter view of full 3D face model (left), mid-sagittal view of tongue, hard palate, and velum 3D models with joint-based rigs (middle), and midsagittal view of face, tongue, hard palate, and velum 3D models with joint-based rigs labeled (right). Key: ul - upper lip, ll - lower lip, li - lower incisor, hp - hard palate, v - velum, t - tongue.

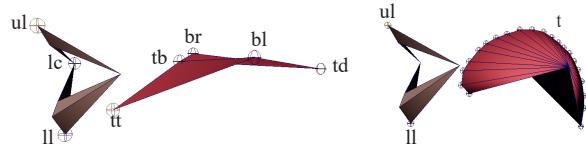


Figure 5. Comparing inner mouth model resolution of the 3D EMA-based model in [33] (left) and the proposed 3D MRI-based model (right). Tongue resolution is increased from 6 joints in [33] to 20 joints in proposed model.

Rather than receiving the full waveform as in the offline case, we instead use an audio input stream to collect batches of 1600 audio samples in sequential order from a WAV file or from a live microphone. Given that our inversion model predicts noisy MRI trajectories from silence, we employ Google’s WebRTC Voice Activity Detector (VAD) to classify if a batch contains speech or not. If not, we do not move the facial model for 100ms as silence simulation.

If the VAD detects speech, we use a sliding context window of n seconds. This window consists of three parts: the *seen* batch, the *current* batch, and the *forward* batch. The current batch contains the 1600 samples that we have just received from the audio input stream. The forward batch is 1600 samples from a look-ahead window of 100ms duration. This audio is “future audio” from an intentional initial 100ms delay, where we wait for two batches to be streamed in before starting speech processing. Using the forward batch, we add future context that helps the inversion model by providing coarticulation information from the next 100ms. Since the current and forward batches together have 3200 samples of speech data, we prepend the $16000n - 3200$ samples of audio that we have just processed as past context in the seen batch. Together, the seen, current, and forward batches make an n -second long context window for the inversion model.

One other issue we encounter is the lack of context when streaming first begins. Since we only have access to a small number of 100ms audio batches as input, the inversion model is prone to output noise. To mitigate these initial fluctuations, we employ a source of artificial context—we prepend either a random utterance from the USC-TIMIT dataset, an articulated vowel, or silence to our streaming audio until our sliding window has a full n seconds of context to draw from.

This context window is translated to the corresponding MRI trajectories using deep articulatory inversion. Of the $83n$ frames returned by the model, we only retain the approximately 8 frames representing the current batch of audio.

6.5. Real-time MRI-to-Avatar

We import the fully rigged 3D facial model designed for vocal visualization into Unreal Engine to fully utilize its optimizations in real-time animation.

To dynamically transform the joints of our model without the use of a Control Rig, we define custom Actor and AnimInstance classes for each MRI segment. Using imported Maya rigs as SkeletalMeshComponents, we can programmatically share our desired feature location with a corresponding Animation Blueprint and control the joint when the game engine is live. Additionally, each joint’s orientation is set relative to the World Space to ensure we preserve covariances between MRI features for physiological accuracy.

To stream the processed MRI trajectory batches to Un-

real Engine, we implement a simple User Datagram Protocol (UDP) between two simultaneous processes. We first perform all speech-to-MRI preprocessing in an async-based Python process then pass information across a socket to the concurrent Unreal Engine process. However, since Ticks in Unreal Engine are constrained to 83.33 Hz to match the MRI frequency, packet loss will severely degrade the animation quality when the server send rate is high. To remedy this timing mismatch, we halve the server send rate to 42 Hz and utilize an asynchronous TQueue within Unreal Engine to collect every frame. We are able to perfectly match the rate of streaming audio in this manner.

With this implementation of real-time MRI-to-avatar in Unreal Engine, we can interact with a low-latency streaming vocal animation in a game-like 3D environment.

7. Results

Examining both the quantitative and qualitative parts of our MRI vocal tract analysis provides valuable insights into both the efficacy of our proposed methods and MRI as a representation for naturalistic speech production. We explore these topics in brief in this section.

7.1. Feature Extraction

When analyzing our various feature extraction methods, we first evaluate performance within the context of seen speakers but unseen examples.

Figure 6 highlights quantitative results in L1 losses and Pearson Correlation Coefficients (PCCs) when evaluating models on unseen examples from seen speakers. We observe a significant trend where multimodal models perform consistently better than the purely video-based U-Net. In fact, the best model in terms of both metrics includes the outputs of the U-Net as one of the input modalities alongside WavLM vectors. These results suggest the inclusion of speech within segmentation provides additional speaker-specific information related to the anatomy of the vocal tract. Since the shape of different parts of the vocal tract can greatly vary from speaker to speaker, this inclusion is crucial to a better in-domain modeling of speech production. With the image modality alone, the fully pixel value-based U-Net generalizes better to unseen speakers since contour pixel values have less dependence on the speaker compared to WavLM features in the speech modality.

For visualization of these results, we invite you to watch our demo video in supplementary materials.

7.2. Deep Speech Synthesis

Similarly, we evaluate our segmentation methods using the MRI-based speech synthesis downstream task within seen and unseen speaker contexts.

To summarize 5.2, the synthesis model is pretrained on the newly-labeled 75-speaker dataset. To then evaluate the

Model	Mean WER [↓]	
	Seen Speaker	Unseen Speaker
U-Net + WavLM	0.31 ± 0.36	0.33 ± 0.26
U-Net	0.36 ± 0.33	0.35 ± 0.33
Ground Truth	0.34 ± 0.35	0.50 ± 0.27
U-Net + WavLM (S)	0.35 ± 0.33	0.50 ± 0.39

Table 1. USC-TIMIT speaker finetuning for seen and unseen speakers: Mean WER for speech synthesis pretrained on 75-speaker dataset. (S) denotes synthesis model pretrained using single MRI speaker. All other models are pretrained with 75-speaker MRI.

performance of a given feature extraction model (e.g. U-Net, multimodal), we finetune this model on the predicted MRI trajectories of a USC-TIMIT speaker.

To evaluate the intelligibility of synthesized speech, we compute the word error rate (WER) on test unseen examples from the same training speaker using Whisper [39], a state-of-the-art automatic speech recognition (ASR) model. For seen speakers of segmentation models, the multimodal UNet-WavLM based synthesizer outperforms both the ground truth baseline as well as the U-Net, suggesting that the addition of the speech modality helps preserves more speech-related information within the predicted MRI point trajectories compared to a purely image-based approach. Table 1 summarizes these results.

Similarly, Table 1 highlights that the UNet-WavLM based model has the lowest WER when testing against an unseen USC-TIMIT speaker, compared to the ground truth segmentations and the U-Net. This demonstrates that the outputs from multi-modal on unseen speaker still capture representative articulatory kinematics for naturalistic speech.

7.3. Visualization

We examine the visualization results in both quantitative and qualitative manners. To evaluate streaming performance, we explore both the latency and the accuracy of the streaming system.

For our baseline, we use the 1-second sliding context window results from [38], where they achieve an average latency of 133ms per 100ms batch of streamed audio. We conduct similar profiling tests with the same 1-second sliding context window and find an average latency of **107ms** per 100ms batch of streamed audio, outperforming the baseline by a 20% margin while having 6× more vertices. We attribute this largely to the animation times being reduced from 56.3ms/batch when streaming in Maya to a constant 12ms/batch using the optimized game engine.

Additionally, we evaluate the streaming performance of

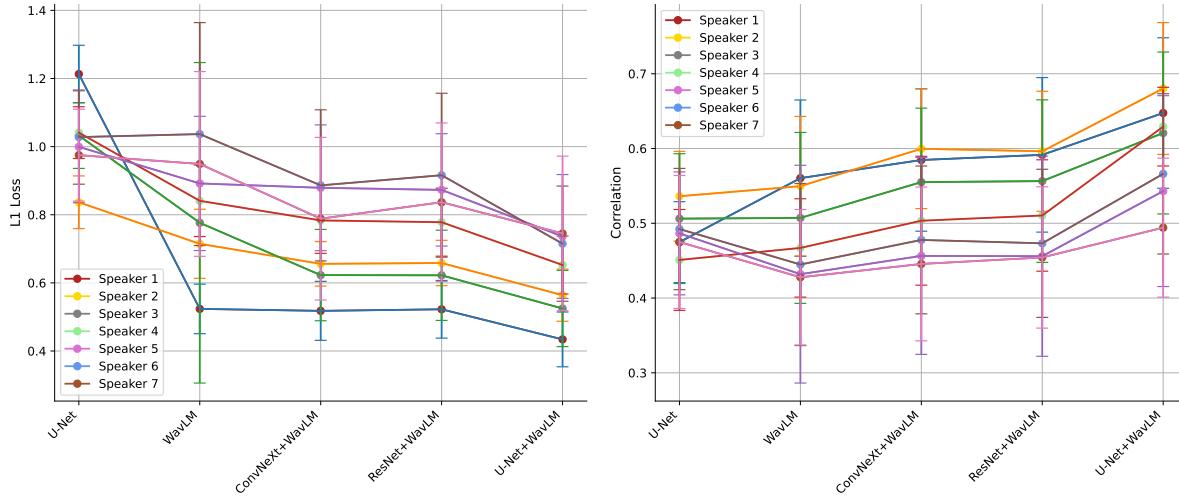


Figure 6. L1 losses [↓] (left) and Pearson Correlation Coefficients (PCCs) [↑] (right) comparing MRI trajectories of unseen examples from seen speakers of a given model with the USC-TIMIT ground truth. Varying through a subset of six representative models.

the inversion model for a quantitative metric on animation accuracy. When streaming an unseen example from a seen USC-TIMIT speaker, Figure 7 highlights the visual similarities between predicted MRI trajectories and predictions from UNet-WavLM. We observe high Pearson correlation

coefficients for the displayed MRI features, which were chosen based on their importance in speech synthesis following [59].

Qualitatively, we visually inspect animation results when streaming and compare these to the ground truth vocal tract movements from the MRI videos. We observe articulator movements matching the U-Net during speech production. For example, Figure 3 demonstrates accurate tongue movement corresponding to the end of the articulation of the word “buns”, where the tip of the tongue first touches the hard palate during the consonant “n” then extends forward during the consonant “s” before receding. We summarize these results in further detail with additional examples in the supplementary video.

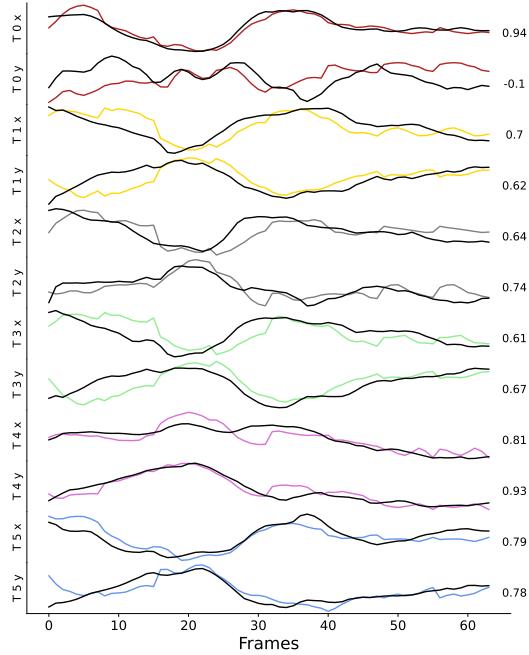


Figure 7. MRI trajectories inferred when streaming the inversion model are shown in color. The trace of the ground truth MRI data is shown in black. To the right of each plot, we present Pearson correlation coefficients (PCCs) comparing the predicted trajectories to the ground truth.

8. Conclusion

Pretraining with audio and image modalities using the newly labeled 75-speaker RT-MRI dataset establishes new MRI benchmarks in vocal tract feature extraction, deep speech synthesis, and real-time speech-driven avatars. While we achieve high quality seen speaker visualization using inversion, current models struggle to disentangle speaker-specific information from speech representations. Future work may use the labeled 75-speaker dataset for speaker-independent speech-driven facial animation.

References

- [1] Anonymized. Towards streaming speech-to-avatar synthesis, 2023. 2, 5
- [2] Anonymized. Deep articulatory mri feature extraction and speech synthesis, 2023. 2, 3

- 592 [3] Anonymized. Acoustic-to-articulatory inversion for multi- 648
593 lingual and downstream tasks, 2023. 2 649
594 [4] Anonymized. Articulatory synthesis with multi-modal and 650
595 self-supervised features, 2023. 4 651
596 [5] Gopala K. Anumanchipalli, Josh Chartier, and Edward F. 652
597 Chang. Speech synthesis from neural decoding of spoken 653
598 sentences. *Nature*, 568(7753):493–498, 2019. 1
599 [6] Monica Villanueva Aylagas, Hector Anadon Leon, Matthias 600
601 Teye, and Konrad Tollmar. Voice2face: Audio-driven facial 602
603 and tongue rig animations with cVAEs. *Computer Graphics Forum*, 41(8):255–265, 2022. 2
604 [7] Narjes Bozorg and Michael T Johnson. Acoustic-to- 605
606 articulatory inversion with deep autoregressive articulatory- 607
608 wavenet. *Networks (CNNs)*, 2020. 2
609 [8] Erik Bresch and Shrikanth Narayanan. Region segmentation 610
611 in the frequency domain applied to upper airway real-time 612
613 magnetic resonance images. *IEEE Transactions on Medical Imaging*, 28(3):323–338, 2009. 1
614 [9] Zexin Cai et al. The dku-jnu-ema electromagnetic articulog- 615
616 raphy database on mandarin and chinese dialects with tandem 617
618 feature based acoustic-to-articulatory inversion. In *ISCSLP*, 2018. 2
619 [10] Claudia Canevari et al. A new italian dataset of parallel 620
621 acoustic and articulatory data. In *Interspeech*, 2015. 2
622 [11] Shicheng Chen, Yifeng Zheng, Chengrui Wu, Guorui Sheng, 623
624 Pierre Roussel, and Bruce Denby. Direct, near real time 625
626 animation of a 3d tongue model using non-invasive ultrasound 627
628 images. In *2018 IEEE International Conference on Acoustics, 629
630 Speech and Signal Processing (ICASSP)*, pages 4994–4998, 2018. 2
631 [12] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, 632
633 Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya 634
635 Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yan- 636
637 min Qian, Yao Qian, Jian Wu, Michael Zeng, Xiangzhan Yu, 638
639 and Furu Wei. Wavlm: Large-scale self-supervised pre- 640
641 training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518, 2022. 2
642 [13] Yaran Chen, Dongbin Zhao, and Haoran Li. Deep kalman 643
644 filter with optical flow for multiple object tracking. In *2019 IEEE International Conference on Systems, Man and Cybernetics (SMC)*, pages 3036–3041, 2019. 3
645 [14] TM Chin, WC Karl, and AS Willsky. Probabilistic and 646
647 sequential computation of optical flow using temporal coherence. *IEEE Trans Image Process*, 1994. 3
648 [15] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag 649
650 Ranjan, and Michael Black. Capture, learning, and synthesis 651
652 of 3D speaking styles. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10101– 653
654 10111, 2019. 2
655 [16] Pif Edwards, Chris Landreth, Eugene Fiume, and Karan 656
657 Singh. Jali: an animator-centric viseme model for expressive 658
659 lip synchronization. *ACM Transactions on Graphics*, 35(4):1–11, 2016. 2
660 [17] Paul Ekman and Wallace V. Friesen. Facial action coding 661
662 system, 1978. 2
663 [18] David Gaddy and Dan Klein. An improved model for voicing 664
665 silent speech. In *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)*, pages 175–181, Online, 666
667 2021. Association for Computational Linguistics. 4
668 [19] Epic Games. <https://www.unrealengine.com/en-US/metahuman-creator>. 2
669 [20] Prasanta Kumar Ghosh et al. A generalized smoothness criterion for acoustic-to-articulatory inversion. *JASA*, 2010. 2
670 [21] Bryan Gick, Barbara May Bernhardt, Penelope Bacsfalvi, 671
672 and Ian Wilson. 11. ultrasound imaging applications in second language acquisition. 2008. 1
673 [22] S Ashwin Hebbar, Rahul Sharma, Krishna Somandepalli, 674
675 Asterios Toutios, and Shrikanth Narayanan. Vocal tract 676
677 articulatory contour detection in real-time magnetic resonance 678
679 images using spatio-temporal context. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7354–7358, 2020. 2, 3
680 [23] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, 681
682 Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman 683
684 Mohamed. Hubert: Self-supervised speech representation 685
686 learning by masked prediction of hidden units. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 29:3451–3460, 2021. 4
687 [24] Thomas Hueber et al. Speaker adaptation of an acoustic-to- 688
689 articulatory inversion model using cascaded gaussian mixture 690
691 regressions. In *Interspeech*, 2013. 2
692 [25] Aravind Illa et al. The impact of cross language on acoustic- 693
694 to-articulatory inversion and its influence on articulatory 695
696 speech synthesis. In *ICASSP*, 2022. 697
698 [26] An Ji. *Speaker independent acoustic-to-articulatory inversion*. PhD thesis, Marquette University, 2014. 2
699 [27] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: 700
701 Generative adversarial networks for efficient and high fidelity 702
703 speech synthesis. In *Advances in Neural Information Processing Systems*, pages 17022–17033. Curran Associates, Inc., 2020. 4
704 [28] Paul K. Krug et al. Self-Supervised Solution to the Control 705
706 Problem of Articulatory Synthesis. In *Interspeech*, 2023. 2
707 [29] June S. Levitt and William F. Katz. The effects of EMA- 708
709 based augmented visual feedback on the English speakers' 710
711 acquisition of the Japanese flap: a perceptual study. In *Proc. Interspeech 2010*, pages 1862–1865, 2010. 1
712 [30] Yongwan Lim, Asterios Toutios, Yannick Bliesener, Ye Tian, 713
714 Sajan Goud Lingala, Colin Vaz, Tanner Sorensen, Miran 715
716 Oh, Sarah Harper, Weiyi Chen, Yoonjeong Lee, Johannes 717
718 Töger, Mairyam Lloréns Monteserin, Caitlin Smith, Bianca 719
720 Godinez, Louis Goldstein, Dani Byrd, Krishna S. Nayak, and 721
722 Shrikanth S. Narayanan. A multispeaker dataset of raw and 723
724 reconstructed speech production real-time mri video and 3d 725
726 volumetric images. *Scientific Data*, 8(1), 2021. 1, 3
727 [31] Peng Liu et al. A deep recurrent approach for acoustic-to- 728
729 articulatory inversion. In *ICASSP*, 2015. 2
730 [32] Yijing Lu, Charlotte E.E. Wiltshire, Kate E. Watkins, Mark 731
732 Chiew, and Louis Goldstein. Characteristics of articulatory 733
734

- 705 gestures in stuttered speech: A case study using real-time
706 magnetic resonance imaging. *Journal of Communication*
707 *Disorders*, 97, 2022. 1, 2
- 708 [33] Salvador Medina, Denis Tome, Carsten Stoll, Mark Tiede,
709 Kevin Munhall, Alex Hauptmann, and Iain Matthews.
710 Speech driven tongue animation. In *2022 IEEE/CVF Conference*
711 *on Computer Vision and Pattern Recognition (CVPR)*.
712 IEEE, 2022. 1, 2, 6
- 713 [34] Sean L. Metzger, Kaylo T. Littlejohn, Alexander B. Silva,
714 David A. Moses, Margaret P. Seaton, Ran Wang, Maximilian E. Dougherty, Jessie R. Liu, Peter Wu, Michael A.
715 Berger, Inga Zhuravleva, Adelyn Tu-Chan, Karunesh Ganguly, Gopala K. Anumanchipalli, and Edward F. Chang.
716 A high-performance neuroprosthesis for speech decoding and
717 avatar control. *Nature*, 620(7976):1037–1046, 2023. 1
- 718 [35] Shrikanth Narayanan, Asterios Toutios, Vikram Rama-
719 narayanan, Adam Lammert, Jangwon Kim, Sungbok Lee,
720 Krishna Nayak, Yoon-Chul Kim, Yinghua Zhu, Louis Gold-
721 stein, Dani Byrd, Erik Bresch, Prasanta Ghosh, Athanasios
722 Katsamanis, and Michael Proctor. Real-time magnetic reso-
723 nance imaging and electromagnetic articulography database
724 for speech production research (tc). *The Journal of the*
725 *Acoustical Society of America*, 136(3):1307–1311, 2014. 1,
726 2
- 727 [36] Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew
728 Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori,
729 Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, Ben
730 Glocker, and Daniel Rueckert. Attention u-net: Learning
731 where to look for the pancreas, 2018. 3
- 732 [37] Slim Ouni et al. Modeling the articulatory space using a
733 hypercube codebook for acoustic-to-articulatory inversion.
734 *JASA*, 2005. 2
- 735 [38] Tejas S. Prabhune, Peter Wu, Bohan Yu, and Gopala K. Anu-
736 manchipalli. Towards streaming speech-to-avatar synthesis,
737 2023. 5, 7
- 738 [39] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman,
739 Christine McLeavy, and Ilya Sutskever. Robust speech
740 recognition via large-scale weak supervision. In *Proceedings*
741 *of the 40th International Conference on Machine Learning*,
742 pages 28492–28518. PMLR, 2023. 7
- 743 [40] Alexander Richard, Colin Lea, Shugao Ma, Jurgen Gall, Fer-
744 nando de la Torre, and Yaser Sheikh. Audio- and gaze-driven
745 facial animation of codec avatars. In *Proceedings of the*
746 *IEEE/CVF Winter Conference on Applications of Computer*
747 *Vision (WACV)*, pages 41–50, 2021. 1
- 748 [41] Korin Richmond. *Estimating articulatory parameters from*
749 *the acoustic speech signal*. PhD thesis, University of Edin-
750 burgh, 2002. 2
- 751 [42] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net:
752 Convolutional networks for biomedical image segmentation,
753 2015. 3
- 754 [43] Kevin Scheck and Tanja Schultz. STE-GAN: Speech-to-
755 Electromyography Signal Conversion using Generative Ad-
756 versarial Networks. In *Interspeech*, 2023. 2
- 757 [44] Endang Setyati, Surya Sumpeno, Mauridhi Hery Purnomo,
758 Koji Mikami, Masanori Kakimoto, and Kunio Kondo.
759 Phoneme-viseme mapping for indonesian language based on
760 blend shape animation. In *IAENG International Journal of*
761 *Computer Science*, 2015. 2
- 762 [45] Abdolreza Sabzi Shahrebabaki et al. Acoustic-to-articulatory
763 mapping with joint optimization of deep speech enhance-
764 ment and articulatory inversion models. *TASLP*, 2021. 2
- 765 [46] Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li,
766 Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei
767 Qin, Jifeng Dai, and Hongsheng Li. Videoflow: Exploiting
768 temporal cues for multi-frame optical flow estimation, 2023.
769 3
- 770 [47] Hayato Shibata et al. Unsupervised acoustic-to-articulatory
771 inversion neural network learning based on deterministic
772 policy gradient. In *SLT*, 2021. 2
- 773 [48] Yashish M Siriwardena et al. The secret source: Incorpor-
774 ating source features to improve acoustic-to-articulatory
775 speech inversion. In *ICASSP*, 2023. 2
- 776 [49] Ingmar Steiner and Slim Ouni. Progress in animation of an
777 ema-controlled tongue model for acoustic-visual speech syn-
778 thesis, 2012. 2
- 779 [50] Atsuo Suemitsu and Jianwu Dang. A real-time articulatory
780 visual feedback approach with target presentation for second
781 language pronunciation learning. *The Journal of the Acous-
782 tical Society of America*, 2015. 1
- 783 [51] Guolun Sun et al. Temporal convolution network based joint
784 optimization of acoustic-to-articulatory inversion. *Applied*
785 *Sciences*, 2021. 2
- 786 [52] Sarah L. Taylor, Moshe Mahler, Barry-John Theobald, and
787 Iain Matthews. Dynamic units of visual speech. In *Pro-
788 ceedings of the ACM SIGGRAPH/Eurographics Symposium*
789 *on Computer Animation*, page 275–284, Goslar, DEU, 2012.
790 Eurographics Association. 2
- 791 [53] Mark Kenneth Tiede et al. Quantifying kinematic aspects of
792 reduction in a contrasting rate production task. *JASA*, 2017.
793 2, 4
- 794 [54] Tomoki Toda et al. Acoustic-to-articulatory inversion map-
795 ping with gaussian mixture model. In *ICSLP*, 2004. 2
- 796 [55] Asterios Toutios and Konstantinos Margaritis. A rough guide
797 to the acoustic-to-articulatory inversion of speech. In *HER-
798 CMA*, 2003. 2
- 799 [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,
800 Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia
801 Polosukhin. Attention is all you need. In *Advances in Neu-
802 ral Information Processing Systems*. Curran Associates, Inc.,
803 2017. 4
- 804 [57] Jianrong Wang et al. Acoustic-to-articulatory inversion
805 based on speech decomposition and auxiliary feature. In
806 *ICASSP*, 2022. 2
- 807 [58] Martijn Wieling et al. Analysis of acoustic-to-articulatory
808 speech inversion across different accents and languages. In
809 *Interspeech*, 2017. 2
- 810 [59] Peter Wu, Tingle Li, Yijing Lu, Yubin Zhang, Jiachen
811 Lian, Alan W Black, Louis Goldstein, Shinji Watanabe, and
812 Gopala K. Anumanchipalli. Deep speech synthesis from mri-
813 based articulatory representations. In *INTERSPEECH 2023*.
814 ISCA, 2023. 1, 2, 3, 4, 8
- 815 [60] Peter Wu et al. Speaker-independent acoustic-to-articulatory
816 speech inversion. In *ICASSP*, 2023. 2

- 819 [61] Xurong Xie et al. Deep neural network based acoustic-to-
820 articulatory inversion using phone sequence information. In
821 *Interspeech*, 2016. 2
- 822 [62] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun,
823 Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven
824 3d facial animation with discrete motion prior. *arXiv preprint*
825 *arXiv:2301.02379*, 2023. 2
- 826 [63] Yuyu Xu, Andrew W. Feng, Stacy Marsella, and Ari Shapiro.
827 A practical and configurable lip sync method for games. In
828 *Proceedings of Motion on Games*, page 131–140, New York,
829 NY, USA, 2013. Association for Computing Machinery. 2
- 830 [64] Tianfang Yan et al. Combining language corpora in
831 a Japanese electromagnetic articulography database for
832 acoustic-to-articulatory inversion. In *Interspeech*, 2023. 2
- 833 [65] Atef Ben Youssef et al. Acoustic-to-articulatory inversion
834 using speech recognition and trajectory formation based on
835 phoneme hidden markov models. In *Interspeech*, 2009. 2
- 836 [66] Yide Yu, Amin Honarmandi Shandiz, and László Tóth. Re-
837 constructing speech from real-time articulatory mri using
838 neural vocoders, 2021. 2, 3