

Distributed Feedforward and Feedback Processing across Perisylvian Cortex Supports Human Speech

Ran Wang^{1*} Xupeng Chen¹

Amirhossein Khalilian-Gourtani¹ Leyao Yu^{2,3}

Patricia Dugan² Daniel Friedman² Werner Doyle⁴

Orrin Devinsky² Yao Wang¹ Adeen Flinker^{2,3*}

Affiliation

¹ Electrical and Computer Engineering Department, New York University, Brooklyn, NY, USA.

² Neurology Department, New York University, New York, NY, USA.

³ Biomedical Engineering Department, New York University, Brooklyn, NY, USA.

⁴ Neurosurgery Department, New York University, New York, NY, USA.

* Corresponding Authors: rw1691@nyu.edu; adeen.flinker@nyu.edu

Abstract

Speech production is a complex human function requiring continuous feedforward commands together with reafferent feedback processing. These processes are carried out by distinct frontal and posterior cortical networks, but the degree and timing of their recruitment and dynamics remain unknown. We present a novel deep learning architecture that translates neural signals recorded directly from cortex to an interpretable representational space that can reconstruct speech. We leverage state-of-the-art learnt decoding networks to disentangle feedforward vs. feedback processing. Unlike prevailing models, we find a mixed cortical architecture in which frontal and temporal networks each process both feedforward and feedback information in tandem. We elucidate the timing of feedforward and feedback related processing by quantifying the derived receptive fields. Our approach provides evidence for a surprisingly mixed cortical architecture of speech circuitry together with decoding advances that have important implications for neural prosthetics.

1 INTRODUCTION

The central sulcus divides the human frontal from the posterior temporal, parietal, and occipital neocortices [34]. Traditionally, this divide separates high order planning and motor execution from sensation. Feedforward execution lies in the frontal cortices in contrast to feedback sensory processing across posterior cortices for the various sensory modalities (e.g., auditory, visual, somatosensory, etc.) [17]. Higher order capacities such as working memory, cognitive control, and decision making are often viewed as initiated by frontal cortices with direct influence on sensory cortices [19, 38, 44].

Human higher order cognitive functions include planning and executing complex speech sequences that carry semantic and linguistic meaning [7, 29]. Speech

28 production is a complex human motor behavior requiring precise coordination of
29 multiple oral, laryngeal and respiratory muscles [42]. These finely tuned motor
30 actions then produce refferent feedback in the auditory, tactile, and proprio-
31 ceptive domains as we process our own speech.

32 The dynamic influence of feedforward commands on sensory feedback is a
33 hallmark of sensory motor systems across the animal kingdom [10]. For exam-
34 ple, motor neurons in cricket both drive the generation of chirping sounds as
35 well as inhibit the auditory system to filter out the loud noise produced by its
36 wings [49]. Similarly, auditory neurons in the marmoset monkey are suppressed
37 during vocalization to provide increased sensitivity to vocal feedback [50]. Pre-
38 vailing models in human speech motor control propose a feedforward system
39 that predicts and generates actions and a feedback system responding to the
40 vocal auditory and somatosensory effects [22, 23, 26–28, 30]. There is a con-
41 sensus that the two systems are anatomically separated, with the feedforward
42 system mainly supported by ventral frontal cortices, while posterior cortices
43 support feedback processing. During feedback processing, frontal cortices need
44 to update new actions dynamically, and it remains unclear which subregions
45 are involved in this process. Moreover, the exact timing of feedforward and
46 feedback engagement across the cortex remains unknown.

47 A growing literature has leveraged unique human electrocorticographic (ECoG)
48 recordings from patients undergoing neurosurgical procedures to obtain a com-
49 bined spatial and temporal resolution critical for investigating speech produc-
50 tion. Studies have detailed the signatures of feedforward speech planning [16]
51 and organization of execution [5, 8] in frontal cortices as well as the subsequent
52 auditory feedback architecture in temporal cortices [15, 20, 21]. To date, evidence
53 of feedback processing has mainly focused on artificially altering the acoustic
54 feedback to create a mismatch in the perceived pitch, providing evidence for en-

55 hanced responses in posterior cortices as well as frontal cortex (i.e., ventral sen-
56 sorimotor cortex) [6]. However, the acoustic perturbation also causes speakers
57 to compensate and change their produced pitch, leading to motor enhancement
58 confounded with the feedback. Recently, the unprecedented signal-to-noise ratio
59 offered by ECoG recordings has ushered deep neural network approaches to de-
60 code speech represented in auditory accurately [1, 3, 46, 47] and sensorimotor [4]
61 cortices. Nevertheless, these approaches have not been able to disentangle feed-
62 forward and feedback contributions during speech production as the motor and
63 sensory responses co-occur.

64 We directly disentangle feedback and feedforward processing during speech
65 production by applying a novel deep learning architecture on human neurosurgi-
66 cal recordings to decode speech (Figure 1). Our approach decodes interpretable
67 speech parameters from cortical signals, which drives a rule-based differentiable
68 speech synthesizer. By learning neural network architectures which apply either
69 casual (predicting using only the past), anticausal (predicting using the future
70 feedback), or both (noncausal), spatial-temporal convolutions, we are able to
71 analyze the overall feedforward and feedback contributions, respectively, as well
72 as to elucidate the temporal receptive fields of recruited cortical regions. In
73 contrast to current models that separate feedback and feedforward cortical net-
74 works, our analyses reveal a surprisingly mixed architecture of feedback and
75 feedforward processing both in frontal and temporal cortices while achieving
76 speech decoding performance on-par or better than previously reported.

77 2 RESULTS

78 We report speech decoding of ECoG data obtained from five participants that
79 took part in a battery of speech production tasks: Auditory Repetition (AR),
80 Auditory Naming (AN), Sentence Completion (SC), Word Reading (WR) and

81 Picture Naming (PN). These were designed to elicit the same set of spoken words
82 across tasks while varying the stimulus modality [41] and provided 50 repeated
83 unique words (400-800 total trials per participant) all of which were analyzed
84 locked to the onset of speech production. We start with an overview of our
85 speech decoding approach.

86 2.1 Speech Decoding Approach

87 • **ECoG Decoder.** The decoder maps the ECoG signals to a set of speech
88 parameters (describing both the voiced and unvoiced components) which
89 are then synthesized to speech spectrograms (Figure 1). The ECoG de-
90 coder architecture is based on recent advances in convolutional neural
91 networks leveraging the ResNet approach [24]. We construct a modified
92 ResNet model with nine layers that treat the cortical input as a spatiotem-
93 poral three-dimensional tensor (two dimensions for the electrode array and
94 one for time, see Methods for details). The decoder is trained such that its
95 output parameters match the reference parameters derived from a speech
96 encoder (which is learnt separately in an unsupervised manner). Further-
97 more, our approach ensures that the speech spectrogram derived from
98 these parameters and constructed by the speech synthesizer matches with
99 the actual speech spectrogram. We use this approach to be more data-
100 efficient and allow us to train on a small set of samples for each patient.

101 • **Speech Parameters.** Our speech representation is motivated by the
102 vocoders used for low-bit-rate speech compression dating back to the
103 1980s. We model speech signals as a mixture of voiced and unvoiced com-
104 ponents, with the voiced component described by a source-filter model
105 (dynamically filtered harmonic signals) [13] and the unvoiced component
106 generated by white noise broadband filtering. In addition to the mixing

107 parameter, our representation includes speech formant information (fre-
108 quency, bandwidth, etc.) and loudness (i.e., the energy of speech). See
109 Methods Figure 6 for details.

110 • **Synthesizer.** We use a set of signal processing equations (such as har-
111 monic oscillation, noise generation, filtering, etc.) to synthesize the spec-
112 trogram from our proposed speech parameters. We can train the ECoG
113 decoder with a limited amount of training data by limiting the number
114 of speech parameters and using differentiable signal processing equations.
115 It is noteworthy that the equations we use are differentiable (see Dif-
116 ferentiable Speech Synthesizer in Extended Data A.1), which allows for
117 backpropagation from the spectrogram to the actual learning of the de-
118 coder.

119 • **Speech Encoder.** The speech encoder is pre-trained using an indepen-
120 dent unsupervised approach before the ECoG decoder training. The en-
121 coder is trained to generate a set of speech parameters from a given spec-
122 trogram, from which the aforementioned speech synthesizer can reproduce
123 the spectrogram. This pre-trained encoder generates reference speech pa-
124 rameters from actual speech signals used for the training of the ECoG
125 decoder. The unsupervised process can be easily used to train the speech
126 encoder from any set of speech signals, including patient-specific speech
127 (see details in the Method section 4.4 and Extended Data A.1). Import-
128 antly, this process constrains the speech parameter space to optimize the
129 training of our ECoG decoder, and the parameters can directly drive a
130 speech synthesizer based on differential equations.

131 We trained three separate models using the proposed pipeline, varying in the
132 causality of the temporal convolution used in the ECoG decoder. The causal
133 model uses only past (up to the current) neural signals to produce the current

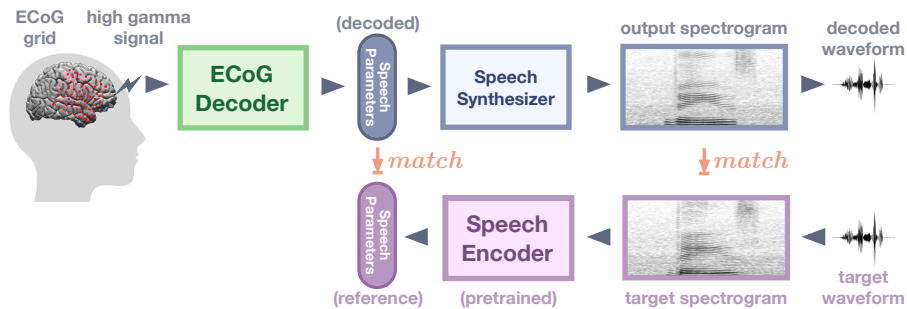


Figure 1: The overall structure of the decoding pipeline. ECoG amplitude signals are extracted in the high gamma range (i.e., 70-150 Hz). The ECoG Decoder translates neural signals from the electrode array to a set of speech parameters. This parameter representation is used to drive a speech synthesizer which creates a spectrogram (and associated waveform). During the training of the ECoG decoder, the speech parameters are matched to a reference derived by a speech encoder pre-trained using an unsupervised approach (without costly manual annotations). This approach constrains the learnt speech parameters and provides naturalistic decoded speeches.

134 speech sample, which reflects feedforward processes. The anticausal model only
135 uses current and future neural signals, reflecting feedback processes. And finally,
136 the non-causal model uses both the past, current, and future neural signals,
137 which are typically used in previous literature and confounds feedforward and
138 feedback processing. The causal and anticausal models allow us to directly assess
139 and tease apart the feedforward and feedback contributions of different cortical
140 regions. It is important to recognize that only causal models are appropriate
141 for real-time speech prosthetic applications.

142 2.2 Speech decoding performance

143 We first demonstrate that our approach produces accurate speech decoding
144 with detailed acoustic features. The model's decoded spectrogram preserves
145 the spectro-temporal structure of the original speech. It reconstructs both vowels,
146 consonants (Figure 2a) as well as the overall spectral energy distribution

147 (Extended Data Figure E1). These acoustic details result in a reconstruction
148 that preserves the speakers' timbre (see Supplementary Video) and leads to
149 naturalistic voice decoding. Our model's speech parameters which include loud-
150 ness, formant frequency, and the mixing parameter (i.e., the relative weighting
151 between voiced and unvoiced components), are decoded accurately with the
152 correct temporal alignment of each word onset and offset (Figure 2b, c). The
153 overall accuracy of the fundamental frequency (i.e., pitch), the first two modeled
154 formants (i.e., F1, F2), and the transition between voiced and unvoiced sounds
155 are a major driving force for accurate speech decoding as well as naturalistic
156 reconstruction that mimics the patient's voice.

157 In order to evaluate the performance and quality of speech, we used sev-
158 eral objective metrics, including the correlation coefficient (CC) between the
159 decoded spectrogram and actual produced speech [2, 3, 25], an objective mea-
160 sure for speech intelligibility known as the Short-Time Objective Intelligibility
161 (STOI) [3, 45], and a measure of spectral distortion, Mel-cepstral distortion
162 (MCD) [4, 35]. Across all participants and metrics, our neural decoding re-
163 sults performed well above chance (Figure 2d in grey; estimated using shuffled
164 data, see Methods section 4.6) and approached an upper bound of performance
165 based on the unsupervised autoencoder (i.e., speech-to-speech) which did not
166 use neural data. The performance range across metrics, and our participants
167 were equal to and often better than the current literature [2–4, 25]. Critically,
168 all these models represent the non-causal case (Figure 2d) that uses data both
169 from the past (feedforward) and the future (feedback), as is currently a common
170 practice [1–4, 37] except a nominal few models [25].

171 In order to directly assess the performance of the causal (predicting using
172 only the past) and anticausal (predicting using the future feedback) models
173 and compare them with the non-causal (using past and future) model, which

174 is standard in the field, we trained three separate models varying the tempo-
175 ral convolution direction. Our results (Figure 2e) show a slight decrease in
176 performance with the causal model. However, it performs close to the other
177 models while providing a causal interpretation, which only uses past signals to
178 predict future speech. This is encouraging, as it suggests that, with additional
179 improvement in the decoder design and training, it is possible to design practi-
180 cally applicable neuroprosthetic speech synthesizers. Also, comparable perfor-
181 mance between causal, anticausal and non-causal approaches indicates a similar
182 amount of information contained by feedforward and feedback signals. Both
183 causal and anticausal models are appropriate for feedforward-feedback analysis
184 and comparison.

185 **2.3 Feedforward and feedback cortical contributions to** 186 **speech production**

187 To elucidate the feedforward and feedback contribution of different cortical re-
188 gions to speech production, we examined the relative contribution of each elec-
189 trode to decoding speech in our models. We derived the relative contribution
190 by quantifying how the input signal at a particular electrode affects the over-
191 all accuracy (measured by the CC) of the reconstructed speech in the causal
192 and anticausal models, respectively (see Methods 4.5). In both the causal and
193 anticausal models, peri-sylvian electrodes were important for speech decoding;
194 however, there was a surprising recruitment of frontal regions when decoding
195 speech based on the feedback (anticausal model, Figure 3b) as well as recruit-
196 ment of temporal sites when decoding speech based on the feedforward signals
197 (causal model, Figure 3c). We only show significant contributions that are
198 above a threshold derived from the shuffled model (depicted in Figure 3d). In
199 order to quantify the prevalence of feedforward or feedback processing, we di-

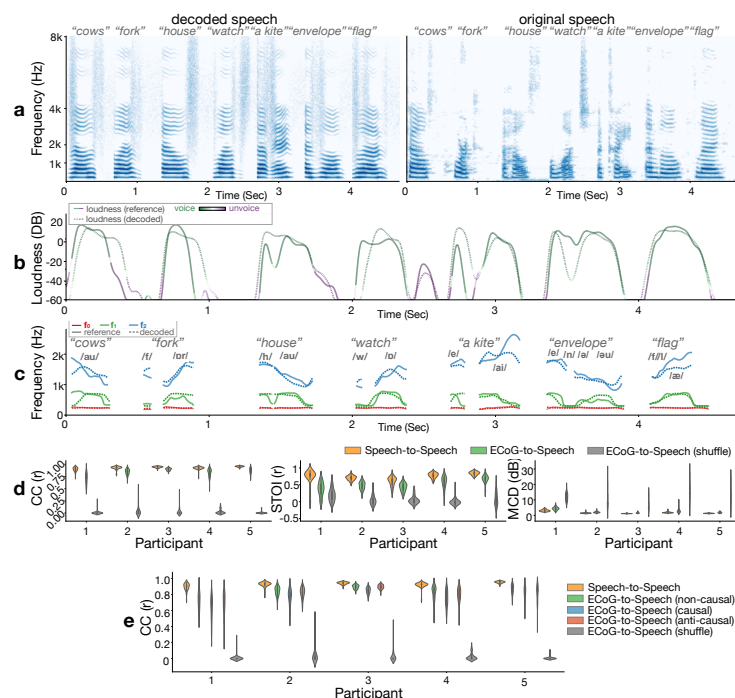


Figure 2: Comparison of original and decoded speech produced by the model. (a) Spectrograms of decoded (left) and original (right) speech exemplar words. (b) Decoded loudness parameter with the voiced (mostly vowel) or unvoiced (mostly consonant) mixing parameter color-coded over the loudness curves. The same color spread and amplitude trend between decoded (dashed) and reference (solid) curves reflect accurate decoding of voice and unvoiced phonemes with correct energy and temporal alignment. (c) Frequencies of the first two formants (F1, F2) and the pitch (F0). The matching between decoded (dashed) curves and reference (solid) curves in both frequencies during each phoneme and the overall temporal dynamic leads to intelligible and naturalistic decoding of voiced sounds. (d) Evaluation of the decoded speech quality in objective metrics. The correlation coefficient of spectrograms (CC, left), short-time objective intelligibility (STOI, middle), and Mel cepstral distortion (MCD, right) are used for the evaluation. Note that lower MCD values represent better performance. Both the reconstructed speech from the speech auto-encoder (yellow) and the speech decoded by the ECoG decoder (green) are reported. Additionally, the performance of a model trained on shuffled data (trained by matching the decoded spectrogram from the neural signal in a given duration to a randomly selected segment of spectrograms during the entire recording session) is also reported as a control. (e) Comparison of the CC metric among noncausal (green), causal (blue), and anticausal (red) models. Compared to the shuffled model (the same shuffled model as in Figure 2d), the comparable performance across noncausal, causal, and anticausal models demonstrates sufficient information for decoding speech from both feedforward and feedback signals during speech production.

200 rectly contrasted the two and projected the results onto the cortex (Figure 3e).
201 To ascertain regions that contribute significantly more to feedback or feedfor-
202 ward processing, we conducted a region of interest analysis, based on within-
203 subject anatomical labels of each electrode (see Methods section 4.3), testing for
204 an increase in causal or anticausal contributions across trials (non-parametric
205 paired Wilcoxon test; Figure 3f). We found a surprisingly mixed distribution of
206 causal and anticausal contributions within both temporal and frontal cortices.
207 A majority of temporal cortex were predominantly anticausal, including caudal
208 superior temporal gyrus (STG; Wilcoxon sign rank, $P=1.607E-15$, $Z=9.6234$)
209 and portions of middle temporal gyrus (MTG; rostral MTG: Wilcoxon sign
210 rank test $P=2.5108E-04$, $Z=4.9359$, and middle MTG: Wilcoxon sign rank test
211 $P=1.5257E-13$, $Z=9.0185$) as well as supramarginal cortex (Wilcoxon sign rank
212 test $P=1.1144E-04$, $Z=5.3919$), implicating it in processing the auditory feed-
213 back signals for speech production. However, there was also a significant causal
214 contribution in rostral STG (Wilcoxon sign rank test $P=0.0332$, $Z=-2.9628$).
215 Similarly, the majority of sensorimotor cortex was predominantly casual, impli-
216 cating it in processing the motor speech commands including ventral precentral
217 (Wilcoxon sign rank, $P=4.9511E-08$, $Z=-7.1409$) and postcentral gyri (Wilcoxon
218 sign rank, $P=6.419E-04$, $Z=-4.9612$). However, the dorsal division of precen-
219 tral gyrus was equally causal and anticausal (Wilcoxon sign rank, $P=0.4349$,
220 $Z=0.6525$), implicating it in processing both feedforward and feedback informa-
221 tion equally. Within the inferior frontal cortex, we found a striking division of
222 function wherein pars opercularis was significantly causal (Wilcoxon sign rank
223 test, $P=8.0693E-15$, $Z=-9.6185$) while pars triangularis was significantly an-
224 ticausal (Wilcoxon sign rank test, $P=2.6715E-06$, $Z=6.3518$). Overall, these
225 findings provide evidence for a mixed feedforward and feedback processing of
226 speech commands and their refference across temporal and frontal cortices, in

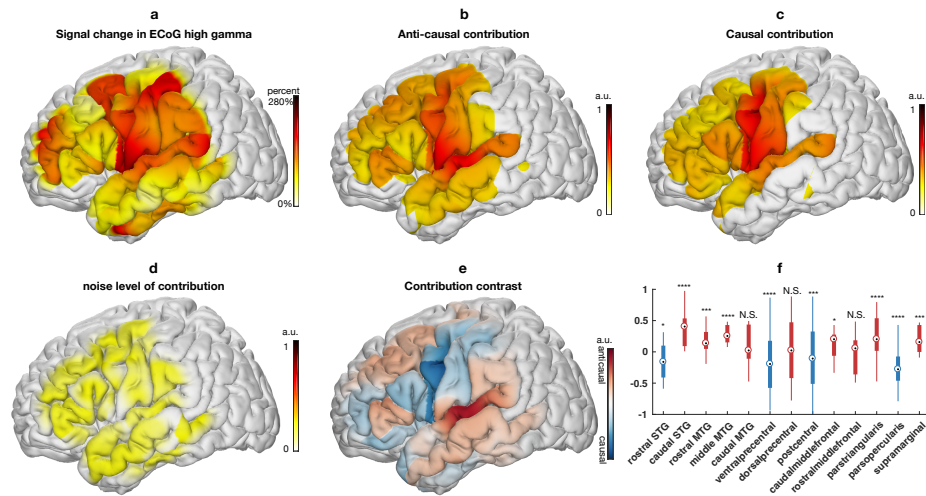


Figure 3: (a) averaged signal of input ECoG projected on the standardized MNI anatomical map. The colors reflect the percentage change of high gamma compared to the baseline level during the pre-stimulus baseline period. (b) shows the anticausal contribution of different cortical locations (red indicates higher contribution), while (c) illustrates the causal contribution. (d) noise level of the contribution analysis evaluated by the contributions from the shuffled model. Contributions below noise level are not shown in (b) and (c). (e) the contrast obtained by taking the difference of the anticausal and causal contribution maps (red means higher anticausal contribution, while blue means higher causal contribution). The boxplots (f) show the average difference in each cortical region (*: P-value<0.05, **: P-value<0.01,***: P-value<0.001,****: P-value<0.0001).

227 contrast to a dichotomous view.

228 2.4 Temporal dynamics and receptive fields of speech pro- 229 duction

230 Speech production includes articulatory planning and executing the motor com-
231 mands, processes that recruit distinct regions of frontal cortex [16]. However,
232 their exact temporal receptive fields remain poorly understood. Earlier, we
233 examined the causal and anticausal cortical contributions during speech artic-
234 ulation. Next, we examine articulatory planning and articulation of speech

235 production stages and derive the related temporal receptive fields across the
236 cortex. We leverage the receptive fields to test how cortical regions contribute
237 differently to speech decoding with time and how frontal cortex dynamics change
238 when feedback is introduced (after articulation starts). Both feedforward and
239 feedback information is processed in tandem.

240 We employed a similar occlusion approach to derive the temporal receptive
241 fields as in the previous section. However, we quantified how the input signal
242 at a particular electrode affects the accuracy of the reconstructed speech across
243 varying delays (see Methods section 4.7). This approach allowed us to quan-
244 tify the contribution of a specific electrode in the model as a function of delay
245 relative to speech decoding, similarly to classical temporal receptive fields (i.e.,
246 TRF). We conducted this analysis for both causal and anticausal models during
247 two epochs – one prior to production (-512ms \sim 0 ms; Figure 4a) and the other
248 during production, which included both causal and anticausal components (0ms
249 \sim 512ms; Figure 4b, c). The projection of all the temporal receptive fields onto
250 the cortex, which were significantly above a threshold derived from the shuffled
251 model, are plotted in Figure 4 as a function of delay. We found an increased
252 frontal and MTG contribution prior to production (Figure 4a) compared with
253 during production (Figure 4b). These processes are likely related to articulatory
254 planning and lexical retrieval prior to speech production. During production,
255 there was a prominent sharpening of ventral precentral gyrus receptive fields
256 marked by a significant increase in contribution compared with pre-production
257 (Wilcoxon sign rank test, $P=8.3979E-05$, $Z=5.4203$). While a majority of pre-
258 frontal regions engaged prior to production, there was a significant decrease in
259 contribution across pars triangularis (Wilcoxon sign rank test, $P=1.8493E-32$,
260 $Z=-13.6074$), middle frontal gyri (MFG; Wilcoxon sign rank test, $P=3.9177E-09$,
261 $Z=-7.6103$ for caudal and $P=4.1581E-04$, $Z=-4.8311$ for rostral) except for pars

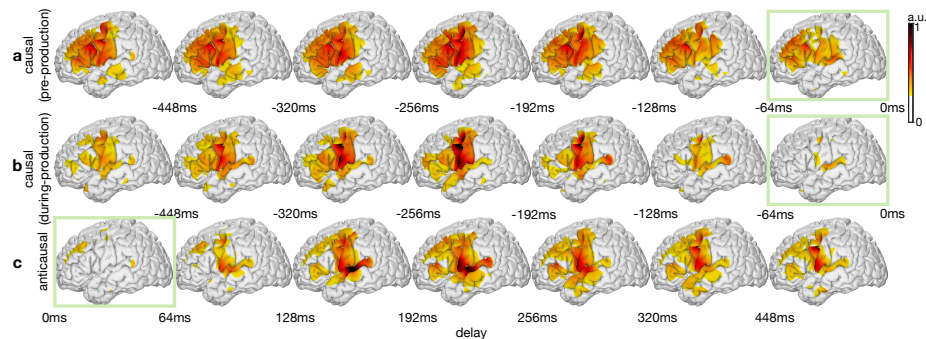


Figure 4: Spatial-temporal receptive fields based on decoding contribution. The contribution to decoding the current speech from cortical neural responses with certain temporal delays. (a) and (b) are the feedforward spatial-temporal receptive fields derived from the causal model by evaluating the contribution of past (negative delays) neural signals during a period before production onset (a) and after onset (b). (c) represents the feedback spatial-temporal receptive fields derived from the anticausal models that evaluate the contribution of future (positive delays) neural signals during feedback after articulation. Contributions below significance ($p < 0.05$) representing the noise level are clipped and not shown in the plots.

262 opercularis (Wilcoxon sign rank test, $P=0.4819$, $Z=0.2066$). Similarly, to our
263 previous results (Figure 3e,f), during production, we found a significant increase
264 in anticausal contribution for caudal STG (Wilcoxon sign rank test, $P=2.6789E-$
265 17 , $Z=9.6711$), pars triangularis (Wilcoxon sign rank test, $P=0.0162$, $Z=3.9003$)
266 and caudal MFG (Wilcoxon sign rank test, $P=0.0045$, $Z=3.9862$) compared with
267 causal contributions. This confirms the anatomical-functional division of the in-
268 ferior and middle frontal gyri as well as caudal (Wilcoxon sign rank test, $P =$
269 $2.6789E-17$, $Z = 9.6711$) and rostral separation of STG (Wilcoxon sign rank
270 test, $p= 0.0343$, $Z= -2.9457$).

271 Next, we conducted a region of interest analysis, based on within-subject
272 anatomical labels of each electrode, in order to derive the temporal receptive
273 curves per region (Figure 5). This approach provides critical insight as to the
274 temporal tuning and peak recruitment of various regions to feedforward process-
275 ing prior to (Figure 5a) and during production (Figure 5b) as well as feedback

276 processing (Figure 5c). We found a shift in receptive field tuning for the two
277 subdivisions of precentral gyrus. Prior to production, dorsal and ventral pre-
278 central gyri were not significantly different from each other (Wilcoxon sign rank
279 test, $P=0.454$, $Z=-0.36103$), and had close peak times (-196ms, -192ms prior
280 to speech for ventral and dorsal precentral gyri, respectively). However, dur-
281 ing production, these dynamics shifted and we found a significant decrease in
282 dorsal precentral causal contribution (Wilcoxon sign rank test, $P=4.7575E-05$,
283 $Z=-5.6272$) accompanied by a temporal separation of peaks (-208ms, -184ms
284 for ventral and dorsal precentral gyri, respectively; Figure 5a,b). Within the
285 inferior frontal gyrus, we found pars opercularis was recruited similarly both
286 prior to production and during production for feedforward processing (Wilcoxon
287 sign rank test, $P=0.5922$, $Z=1.7462$) at a peak delay of -248ms and -280ms,
288 respectively. During production, pars triangularis had a selective increase in
289 recruitment for anticausal compared with causal contributions (Wilcoxon sign
290 rank test, $P=0.0162$, $Z=3.9003$), implicating it in increased feedback processing
291 (Figure 4c, Extended Data Tables 2, 3). The anticausal receptive fields during
292 production provide evidence for feedback processing most strongly contributed
293 by caudal STG, with the earliest peak in contributions seen in dorsal precentral
294 gyrus (144 ms) and caudal STG (168 ms) followed by parietal (supramarginal
295 184ms, postcentral 192ms) and ventral precentral (280 ms) gyri (Extended Data
296 Table 3). These findings suggest a preferential recruitment of prefrontal cortices
297 in feedforward processing prior to production followed by a shift in dynamics
298 during production when feedforward and feedback signals are jointly processed
299 with anatomical divisions of labor.

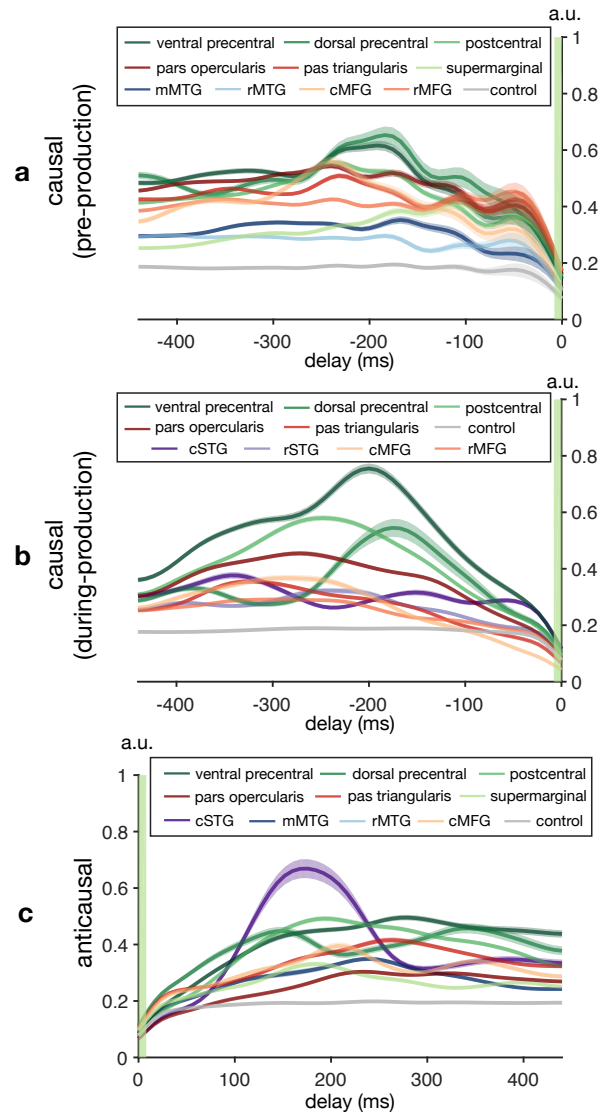


Figure 5: The temporal receptive field across anatomical regions. The contribution to decoding the current speech from cortical neural responses with certain temporal delays. (a) and (b) are the feedforward temporal receptive fields derived from the causal model by evaluating the contribution of past (negative delays) neural signals during a period before production onset (a) and after onset (b). (c) represents the feedback temporal receptive fields derived from the anticausal models that evaluate the contribution of future (positive delays) neural signals during feedback after articulation. The temporal propagation of the shuffled model estimates the noise level dynamics (grey curves in plots). Only regions significantly above noise level (Wilcoxon sign rank test on across-time averaged data, $P < 0.05$) are reported.

300 **3 DISCUSSION**

301 Our study leverages a novel deep learning approach together with neurosurgical
302 recordings and, to our knowledge, is the first to dissociate direct feedforward and
303 feedback cortical contributions during speech production. Our neural network
304 architecture achieves state-of-the-art decoding of speech production, by tapping
305 an interpretable compact speech representation and can be altered to focus on
306 causal, anticausal and non-causal decoding. Our analyses of the cortical contri-
307 butions driving the performance of these models reveal a mixed distribution of
308 feedforward and feedback processing during speech production. This was promi-
309 nent in inferior, middle frontal, and superior temporal gyri which exhibited an
310 anatomical division between feedforward and feedback processing. Lastly, we
311 show a change in the temporal dynamics of frontal recruitment during speech
312 planning through production, characterized by a shift of inferior frontal and pre-
313 central gyri recruitment, processing both feedforward and feedback information
314 at different time points and spatial locations.

315 A growing number of studies have leveraged deep neural networks for cor-
316 tical speech decoding. Convolutional neural networks (CNN) [1, 3, 46, 47] and
317 recurrent neural networks (RNN) [4] have mapped ECoG signals into speech and
318 text [37]. However, our approach diverges from these studies. *Firstly*, we develop
319 a novel differentiable speech synthesizer that can generate natural speech from a
320 compact set of interpretable speech parameters based on several signal process-
321 ing equations. This rule-based synthesizer allows for unsupervised pre-training
322 of meaningful encoded representations (reference speech parameters), as well as
323 reduces the capacity of the entire model and increases training data efficiency.
324 Our approach provides a direct mapping to a patient’s voice. It eliminates the
325 need for labeled articulatory data that maps speech to articulatory dynamics as
326 proposed by Anumanchipalli et al. [4]. *Secondly*, our compact speech representa-

327 tion leverages an interpretable decomposition of speech into voiced and unvoiced
328 components. This decomposition is biologically necessary, has been reported in
329 neural representations across frontal and temporal cortices [8, 31] and stands
330 in contrast to other traditional speech synthesizing approaches [13, 14]. *Lastly*,
331 the speech neural decoding models to date mostly employ non-causal operations.
332 Since such decoders require both past and future information for decoding, they
333 are not applicable for real-time speech prosthetic application. Further, mixed
334 operations hinder disentangling feedforward and feedback cortical contributions.
335 In addition to providing a causal model which directly translates to practical
336 speech prosthetics, our approach provides one of the first reports that can disso-
337 ciate feedforward and feedback cortical contributions during speech production.

338 During speech production, we process feedforward and feedback signals in
339 tandem. It was previously impossible to disentangle the two. Attempts have
340 focused on experimental manipulations which change the feedback by shifting
341 frequency [6] or time [39]. However, these manipulations change the cortical dy-
342 namics and introduce other cognitive processes due to hearing one’s own voice
343 altered as well as induced motor compensation. We applied convolution filters
344 with different causality to directly train models to disentangle feedforward (i.e.,
345 causal models) and feedback (i.e., anticausal models) contributions of cortical
346 regions. Feedforward and feedback processes are critical for driving articulatory
347 vocal tract movement. The feedforward pathway generates an initial articula-
348 tory command and predicts sensory (auditory and somatosensory) targets; the
349 feedback pathway compares the targets with the perceived sensory feedback
350 and updates subsequent feedforward commands. The exact mapping between
351 anatomical regions and their contribution to specific functional roles differ across
352 speech motor control models ([23], [30]). Further, these findings have been de-
353 veloped based mostly on non-invasive studies which have low temporal (e.g.,

354 fMRI) or spatial resolution (e.g., M/EEG). Our high spatio-temporal resolution
355 ECoG data together with advanced deep neural networks provides a fine-grained
356 mapping of the cortical feedforward and feedback speech networks.

357 Consistent with the predominant speech motor control models, our results
358 showed a dominant feedforward process in the ventral motor and pars opercu-
359 laris of the inferior frontal gyrus, while posterior superior temporal and supra-
360 marginal gyri in the parietal lobe showed feedback. However, in contrast to these
361 models, we found that cortices in the frontal lobe, including pars triangularis
362 and caudal middle frontal, are predominantly feedback in nature, while ro-
363stral STG appears feedforward. This feedback processing across frontal cortices
364 became even stronger when we limited our analyses to the speech production
365 epoch (Figure 4c, Extended Data Table 3). Additionally, most gyri (inferior
366 frontal, caudal middle frontal, superior temporal, precentral, and postcentral
367 cortices, see Extended Data Table 2) had both feedforward and feedback con-
368 tributions above the noise level derived from the shuffled model, suggesting the
369 feedforward and feedback processing can mix in these regions.

370 Our results highlight the anticausal feedback signature exhibited by senso-
371 rimotor and frontal cortices. While this goes against the canonical model of
372 the frontal cortex in an action-perception loop [18], our findings complement
373 a growing body of evidence showing specific responses in the frontal cortex to
374 auditory stimuli during perception. Cheung et al. [9] found distinct auditory
375 receptive fields as well as robust passive listening responses in ventral precentral
376 gyrus. Similarly, the dorsal division of precentral gyrus has recently been im-
377 plicated in processing auditory feedback of altered speech as well as responding
378 robustly during passive listening [39]. However, this begs the question as to why
379 the speech motor cortex is processing auditory information. Our feedback con-
380 tribution analysis suggests that the auditory processing is specifically leveraged

381 for anticausal processing of the refferent signals during production. Indeed, our
382 results show that dorsal precentral gyrus decreases feedforward processing while
383 engaged in actual speech production (Figure 5b) and is recruited for feedback
384 at an early time point together with temporal cortices (Figure 5c). Under this
385 view, the auditory frontal responses seen during passive listening may constitute
386 a representation dedicated to feedback processing when speech is produced.

387 To summarize, we provided a new approach to decode speech production and
388 interrogate the recalcitrant problem of mixed feedforward and feedback process-
389 ing during speech production. We were able to leverage feedforward processing
390 only in causal models to drive neural speech prosthetics (as opposed to the lit-
391 erature using non-causal processing [1–4, 37]) as well as provide insights into
392 the underpinning cortical drivers. Our results suggest a mixed cortical archi-
393 tecture in frontal and temporal cortices that dynamically shifts and processes
394 both feedforward and feedback signals across the cortex in contrast to previous
395 views associating feedforward or feedback processing of speech with primarily
396 anterior and posterior cortices, respectively.

397 4 METHODS

398 4.1 Participants and experiments

399 The brain activity data were obtained from five patients, including three fe-
400 male and two male native English speakers, undergoing treatment for refrac-
401 tory Epilepsy at NYU Langone hospital, with implanted electrocorticographic
402 (ECoG) subdural electrode grids. All experimental procedures were approved
403 by the NYU Grossman School of Medicine Institutional Review Board. Pa-
404 tients were provided written and oral consent at least one week prior to surgery
405 by a research team member after separate consultation with the clinical care

406 provider. The subjects were instructed to complete five tasks to pronounce the
407 target words in response to certain auditory or visual stimuli. The five tasks
408 were:

- 409 • Auditory Repetition (AR, i.e., to repeat the auditory words).
- 410 • Auditory Naming (AN, i.e., name a word based on an auditory presented
411 definition sentence).
- 412 • Sentence Completion (SC, i.e., complete the last word of an auditorily
413 presented sentence).
- 414 • Visual Reading (VR, i.e., read aloud visually presented word in written
415 form).
- 416 • Picture Naming (PN, i.e., naming a word based on a visually presented
417 color line drawing).

418 Each task contained the same 50 unique target words while varying stimulus
419 modalities (auditory, visual, etc.). Each word appeared once in the AN and SC
420 tasks and twice in the other tasks. For Participants 1-3, the five tasks included
421 400 trials of the produced words and the corresponding ECoG recordings. The
422 produced speech in each trial has an average duration of 500 ms. We repeated
423 the same five tasks twice for Participant 4 and collected data from 800 trials.
424 For Participant 5, we collected 800 trials by repeating the tasks twice, and we
425 also ran an additional AR task (200 trials) which provided 1000 trials in total.

426 **4.2 Data collection and preprocessing**

427 A microphone recorded the subject's speech during the tasks and was synchro-
428 nized to the clinical Neuroworks Quantum Amplifier (Natus Biomedical, Ap-
429 pleton, WI), which records ECoG. The recordings sampled peri-sylvian cortex,

430 including STG, IFG, pre-central, and postcentral gyri. The ECoG implanted
431 array included standard 64 clinical 8×8 macro contacts (10 mm spacing) as well
432 as 64 additional interspersed smaller electrodes (1 mm) between the macro con-
433 tacts (providing 10 mm center-to-center spacing between macro contacts and 5
434 mm center-to-center spacing between micro/macro contacts, PMT corporation,
435 Chanassen, MN). This FDA-approved array was manufactured for the study,
436 and a member of the research team explained to patients that the additional
437 contacts were for research purposes during consent. The ECoG arrays were im-
438 planted on the left hemisphere in all participants' brains and placement location
439 was solely dictated by clinical care. We trained separate sets of decoding mod-
440 els for each participant. We randomly selected 50 out of all trials from the five
441 tasks for testing and used the remaining data for training. The results reported
442 are for testing data.

443 Each electrode sampled ECoG signals at 2048 Hz, which was decimated
444 to 512 Hz prior to processing. After rejecting electrodes with artifacts (i.e.,
445 line noise, poor contact with cortex, and high amplitude shifts), we subtracted
446 a common average reference (across all valid electrodes and time) from each
447 individual electrode. Electrodes with inter-ictal and epileptiform activity were
448 removed from the analysis (note that the large number of temporal electrodes
449 were removed from patients 4 and 5 for this reason). We then extracted the
450 envelope of the high gamma (70-150 Hz) component from the raw signal with
451 the Hilbert transform and further downsampled it to 125 Hz. The signal of
452 each electrode over the silent baseline of 250 ms before the stimulus was used as
453 the reference signal, and each electrode's signal was normalized to the reference
454 mean and variance (i.e., z-score).

455 **4.3 Electrode localization**

456 Electrode localization in subject space, as well as MNI space, was based on
457 coregistering a preoperative (no electrodes) and postoperative (with electrodes)
458 structural MRI (in some cases, a postoperative CT was employed depending
459 on clinical requirements) using a rigid-body transformation. Electrodes were
460 then projected to the surface of the cortex (preoperative segmented surface) to
461 correct for edema-induced shifts following previous procedures [48] (registration
462 to MNI space was based on a non-linear DARTEL algorithm). Based on the
463 subject's preoperative MRI, the automated FreeSurfer segmentation (Destrieux)
464 is used for labeling within subject anatomical locations of electrodes.

465 **4.4 Speech decoding framework**

466 The backbone of our neural decoding framework is constructed by an ECoG
467 decoder and a speech synthesizer (Figure 6a or Figure 1). During testing, from
468 the high gamma components of the ECoG signal, the decoder generates a set
469 of speech parameters that drive a differentiable speech synthesizer to gener-
470 ate speech spectrograms (and corresponding waveforms by the Griffin-Lim al-
471 gorithm). Besides being trained to work with the speech synthesizer to out-
472 put spectrograms matching the target spectrograms, the ECoG decoder is also
473 trained to match its output with a set of reference speech parameters. This ref-
474 erence matching training strategy provides a more direct gradient to the ECoG
475 decoder such that it converges faster and is less prone to overfitting.

476 The reference speech parameters are derived from a pre-trained speech en-
477 coder. During pre-training, the speech encoder and the speech synthesizer ful-
478 fill an auto-encoding task (i.e., mapping the input spectrogram to the speech
479 parameters and back to the spectrogram) (Figure 6b). When such speech-
480 to-speech reconstruction is accurate, the parameters generated by the speech

481 encoder should provide physically meaningful speech parameters. Since the
482 pre-training is unsupervised and the subject speech audio data is easy to col-
483 lect, obtaining the reference speech parameters is straightforward. Note that
484 the speech-to-speech autoencoder and the reference parameters are only used
485 for the training of the ECoG decoder. Once the ECoG decoder is trained, the
486 trained decoder and the speech synthesizer can be used to convert ECoG signals
487 to speech without the need for reference parameters.

488 More details of the structure of the speech synthesizer (Figure 6e), ECoG
489 decoder (Figure 6c), Speech encoder (Figure 6d), and loss can be found in
490 Extended Data A.1.3.

491 **4.5 Revealing delay-dependent contribution of different** 492 **cortical regions from the trained ECoG to speech model**

Before formally defining the various contribution scores, we introduce the fol-
lowing notations: $A_{\text{ref}}[s]$: the reference spectrogram over a time duration S
centered at time s , i.e., from $s - S/2$ to $s + S/2$, derived by the speech-to-speech
autoencoder. $A_{\text{intact}}[s]$: the model output with “intact” input (i.e., all ECoG
signals are used). $A_{\text{occlude}}^i[s|t]$: the model output at time duration centered at
 s when the i th ECoG electrode signal in the time duration centered at t from
 $t - \frac{T}{2}$ to $t + \frac{T}{2}$ is occluded. $r(\cdot, \cdot)$: correlation coefficient between two signals.
We define the contribution of i th electrode in time duration centered at t to the
output over duration centered at s by the reduction in the correlation coefficient
between the output signal with the reference signal over the duration s when
the i th electrode signal in duration t is occluded. Specifically:

$$C^i[s, t] = \text{Mean}\{r(A_{\text{ref}}[s], A_{\text{intact}}[s]) - r(A_{\text{ref}}[s], A_{\text{occlude}}^i[s|t])\}$$

493 where $\text{Mean}\{\cdot\}$ denotes averaging across all testing samples.

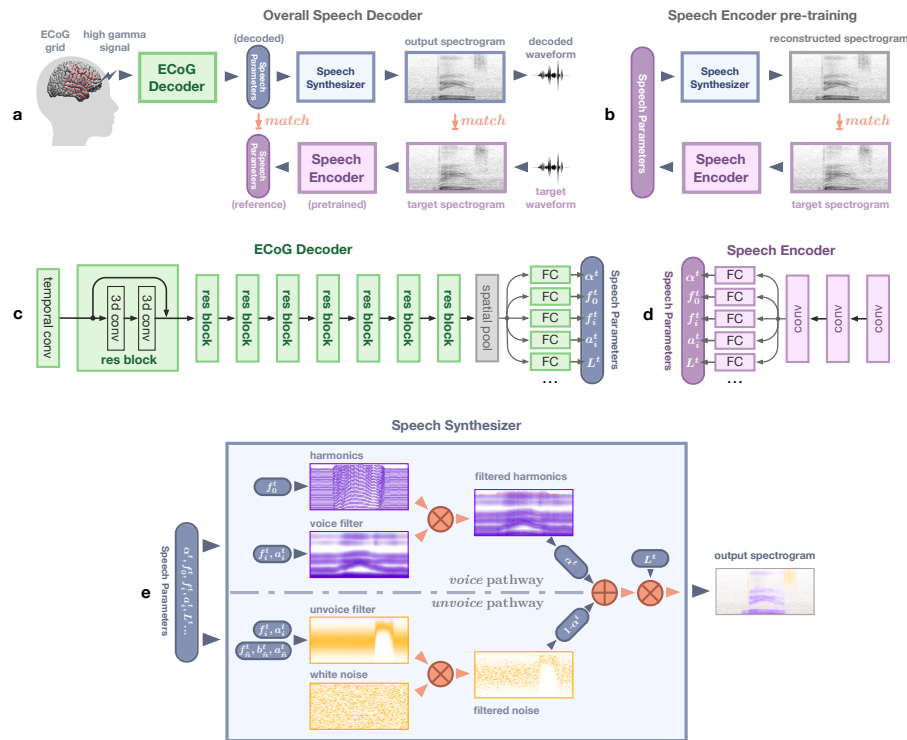


Figure 6: Structure of the decoding framework. (a) The overall network architecture (same as in Figure 1, repeated here for ease of understanding of the architecture). (b) The auto-encoder used to pretrain the speech encoder. The speech encoder is trained to generate proper speech parameters that can reconstruct input spectrograms through the speech synthesizer. (c) The ECoG decoder is a modified spatio-temporal residual network. After an initial temporal convolutional layer and eight residual blocks (constructed by three-dimensional convolution layers), multiple subnetworks (using one or two fully connected layers) generate speech parameters separately. (d) The speech encoder in (b) has three convolutional layers followed by the same multi-head output structure as in (c). (e) Illustrates the processes within the speech synthesizer. The harmonics (in voice pathway) and white noise (unvoice pathway) are generated and filtered (multiplication in spectrogram domain) by voice and unvoice filters, respectively. The filtered results are then weighted averaged according to the mixing parameter and then amplified by the loudness parameter.

494 To generate the contribution map, we first determine the contribution of
495 each electrode (with a corresponding location in the MNI coordinate), which is
496 then diffused into the surrounding area in the same anatomical region using a
497 Gaussian kernel. Since our ECoG grid has hybrid density, to remove the effect
498 of non-uniform density on the diffused result, we normalize the result of each
499 region by the local grid density. The results shown in Figures 3,4, and 5 are
500 obtained by averaging the contribution maps obtained for all test samples for
501 all participants.

502 4.6 Visualizing spatial contribution map

503 The contribution of the entire signal at the i -th electrode to the entire output
504 signal, C^i , is obtained by using the method in Section 4.5 with S and T cov-
505 ering the entire input and output signal duration. The causal and anticausal
506 contribution plots in Figure 3 are generated by applying such analysis to the
507 learned anticausal model (Figure 3b) and causal model (Figure 3c), respectively.
508 The contrast of the anticausal and causal contribution (Figure 3e) for each is
509 the difference between the causal and anticausal contribution map. The noise
510 level for the contribution analysis (Figure 3d) is generated from the shuffled
511 model using non-causal processing (the shuffled model is trained on an artificial
512 dataset with temporal misaligned input-output, and hence models of different
513 causality are equivalent). To generate per region feedback-feedforward box plot
514 (Figure 3f), we calculate the contrast contributions averaged over electrodes of
515 the same within-subject anatomical labels corresponding to each region.

The contrast of the anticausal and causal contribution (as is shown in Figure 3e) of electrode i is defined as

$$\tilde{C}_{contrast}^i = \tilde{C}_{anticausal}^i - \tilde{C}_{causal}^i$$

In order to examine electrode polarization to anticausal or causal contribution, we calculate the normalized version of anticausal and causal contribution contrast:

$$\tilde{C}_{polar}^i = \frac{\tilde{C}_{anticausal}^i - \tilde{C}_{causal}^i}{\tilde{C}_{anticausal}^i + \tilde{C}_{causal}^i}$$

516 By normalizing the contrast of anticausal and causal contribution, \tilde{C}_{polar}^i em-
517 phasize the angle of contribution directing towards anticausal or causal, rather
518 than their contrast. This is what is visualized in Figure E2 a,b in Extended
519 Data (only for those electrodes with either anticausal contribution attribute
520 ($\tilde{C}_{anticausal}^i$) or causal contribution attribute (\tilde{C}_{causal}^i) above noise level deter-
521 mined by the shuffled model). This is what is shown in Extended Data Figure
522 E2.

523 4.7 Visualizing spatial-temporal contribution receptive field

When evaluating the contribution over a finite duration we use small temporal scope $S = T = 64\text{ms}$. To Evaluate the contribution of an electrode signal to the output with various delay, denoted by τ , we average $C^i[s, s + \tau]$ for all s in a certain duration leading to

$$\tilde{C}^i(\tau) = \frac{1}{s_1 - s_0} \sum_{s=s_0}^{s_1} C^i[s, s + \tau]$$

524 Here we assume the effect of delay is independent of actual output time s .
525 When $\tau \leq 0$, $\tilde{C}_{causal}^i(\tau)$ reveals the causal contribution of electrode i to the
526 output (Figure 4 a,b). To investigate pre-production contribution, we restrict
527 $s + \tau$ and s to be no later than the onset of production (vise-versa for during-
528 production analysis). When $\tau \geq 0$ the $\tilde{C}_{anticausal}^i(\tau)$ reveals the anticausal
529 contribution (Figure 4c). This is how the results in Figure 4 were generated,
530 where the causal (resp. anticausal) contribution is derived from the causal (resp.

531 anticausal) model.

532 **4.7.1 Visualizing per region temporal contribution receptive field**

533 Similar to the per region plot in Figure 3f, to generate a temporal contribution
534 curve for each region (Figure 5), we average the spatial-temporal receptive field
535 data (Figure 4) over to the same within-subject anatomical region labels. The
536 control curve is generated by applying the same method for the shuffled model
537 (grey curves in Figure 4). We omit those curves that are not significantly above
538 noise level by Wilcoxon sign rank testing between averaged (over time) region
539 contribution curves and the averaged (over time) noise level curve (see Extended
540 Data Table 2).

541 **References**

- 542 [1] Hassan Akbari, Bahar Khalighinejad, Jose L Herrero, Ashesh D Mehta,
543 and Nima Mesgarani. Towards reconstructing intelligible speech from the
544 human auditory cortex. Scientific reports, 9(1):874, 2019.
- 545 [2] Miguel Angrick, Christian Herff, Garrett Johnson, Jerry Shih, Dean
546 Krusienski, and Tanja Schultz. Interpretation of convolutional neural
547 networks for speech spectrogram regression from intracranial recordings.
548 Neurocomputing, 342:145–151, 2019.
- 549 [3] Miguel Angrick, Christian Herff, Emily Mugler, Matthew C Tate, Marc W
550 Slutzky, Dean J Krusienski, and Tanja Schultz. Speech synthesis from
551 ecog using densely connected 3d convolutional neural networks. Journal of
552 neural engineering, 16(3):036019, 2019.

- 553 [4] Gopala K Anumanchipalli, Josh Chartier, and Edward F Chang. Speech
554 synthesis from neural decoding of spoken sentences. Nature, 568(7753):493–
555 498, 2019.
- 556 [5] Kristofer E Bouchard, Nima Mesgarani, Keith Johnson, and Edward F
557 Chang. Functional organization of human sensorimotor cortex for speech
558 articulation. Nature, 495(7441):327–332, 2013.
- 559 [6] Edward F Chang, Caroline A Niziolek, Robert T Knight, Srikantan S Na-
560 garajan, and John F Houde. Human cortical sensorimotor network under-
561 lying feedback control of vocal pitch. Proceedings of the National Academy
562 of Sciences, 110(7):2653–2658, 2013.
- 563 [7] Edward F Chang, Kunal P Raygor, and Mitchel S Berger. Contemporary
564 model of language organization: an overview for neurosurgeons. Journal of
565 neurosurgery, 122(2):250–261, 2015.
- 566 [8] Josh Chartier, Gopala K Anumanchipalli, Keith Johnson, and Edward F
567 Chang. Encoding of articulatory kinematic trajectories in human speech
568 sensorimotor cortex. Neuron, 98(5):1042–1054, 2018.
- 569 [9] Connie Cheung, Liberty S Hamilton, Keith Johnson, and Edward F Chang.
570 The auditory representation of speech sounds in human motor cortex. Elife,
571 5:e12577, 2016.
- 572 [10] Trinity B Crapse and Marc A Sommer. Corollary discharge across the
573 animal kingdom. Nature Reviews Neuroscience, 9(8):587–600, 2008.
- 574 [11] Li Deng and Douglas O’Shaughnessy. Speech processing: a dynamic and
575 optimization-oriented approach. CRC Press, 2018.

- 576 [12] Jesse Engel, Lamtharn Hantrakul, Chenjie Gu, and Adam Roberts. DDSP:
577 Differentiable digital signal processing. arXiv preprint arXiv:2001.04643,
578 2020.
- 579 [13] James L Flanagan. Speech analysis synthesis and perception, volume 3.
580 Springer Science & Business Media, 2013.
- 581 [14] Mario Fleischer, Silke Pinkert, Willy Mattheus, Alexander Mainka, and
582 Dirk Mürbe. Formant frequencies and bandwidths of the vocal tract trans-
583 fer function are affected by the mechanical impedance of the vocal tract
584 wall. Biomechanics and modeling in mechanobiology, 14(4):719–733, 2015.
- 585 [15] Adeen Flinker, Edward F Chang, Heidi E Kirsch, Nicholas M Barbaro,
586 Nathan E Crone, and Robert T Knight. Single-trial speech suppression of
587 auditory cortex activity in humans. Journal of Neuroscience, 30(49):16643–
588 16650, 2010.
- 589 [16] Adeen Flinker, Anna Korzeniewska, Avgusta Y Shestyuk, Piotr J
590 Franaszczuk, Nina F Dronkers, Robert T Knight, and Nathan E Crone.
591 Redefining the role of broca’s area in speech. Proceedings of the National
592 Academy of Sciences, 112(9):2871–2875, 2015.
- 593 [17] Joaquin M Fuster. The prefrontal cortex—an update: time is of the essence.
594 Neuron, 30(2):319–333, 2001.
- 595 [18] Joaquin M Fuster. Upper processing stages of the perception–action cycle.
596 Trends in cognitive sciences, 8(4):143–145, 2004.
- 597 [19] Joaquín M Fuster. The prefrontal cortex in the neurology clinic. Handbook
598 of clinical neurology, 163:3–15, 2019.
- 599 [20] Jeremy DW Greenlee, Roozbeh Behroozmand, Charles R Larson, Adam W
600 Jackson, Fangxiang Chen, Daniel R Hansen, Hiroyuki Oya, Hiroto

- 601 Kawasaki, and Matthew A Howard III. Sensory-motor interactions for
602 vocal pitch monitoring in non-primary human auditory cortex. PloS one,
603 8(4):e60783, 2013.
- 604 [21] Jeremy DW Greenlee, Adam W Jackson, Fangxiang Chen, Charles R
605 Larson, Hiroyuki Oya, Hiroto Kawasaki, Haiming Chen, and Matthew A
606 Howard III. Human auditory cortical activation during self-vocalization.
607 PloS one, 6(3):e14744, 2011.
- 608 [22] Frank H Guenther. A neural network model of speech acquisition and motor
609 equivalent speech production. Biological cybernetics, 72(1):43–53, 1994.
- 610 [23] Frank H Guenther. Neural control of speech. Mit Press, 2016.
- 611 [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual
612 learning for image recognition. In Proceedings of the IEEE conference on
613 computer vision and pattern recognition, pages 770–778, 2016.
- 614 [25] Christian Herff, Lorenz Diener, Miguel Angrick, Emily Mugler, Matthew C
615 Tate, Matthew A Goldrick, Dean J Krusienski, Marc W Slutzky, and
616 Tanja Schultz. Generating natural, intelligible speech from brain activity
617 in motor, premotor, and inferior frontal cortices. Frontiers in neuroscience,
618 13:1267, 2019.
- 619 [26] Gregory Hickok. Computational neuroanatomy of speech production.
620 Nature reviews neuroscience, 13(2):135–145, 2012.
- 621 [27] Gregory Hickok. The cortical organization of speech processing: Feed-
622 back control and predictive coding the context of a dual-stream model.
623 Journal of Communication Disorders, 45(6):393–402, 2012. 21st Annual
624 NIDCD-Sponsored ASHA Research Symposium (2011):Neuroplasticity in
625 the Mature Brain.

- 626 [28] Gregory Hickok. The architecture of speech production and the role of
627 the phoneme in speech processing. Language, Cognition and Neuroscience,
628 29(1):2–20, 2014.
- 629 [29] Gregory Hickok and David Poeppel. The cortical organization of speech
630 processing. Nature Reviews Neuroscience, 8(5):393, 2007.
- 631 [30] John F Houde and Srikantan S Nagarajan. Speech production as state
632 feedback control. Frontiers in human neuroscience, 5:82, 2011.
- 633 [31] Colin Humphries, Merav Sabri, Kimberly Lewis, and Einat Liebenthal.
634 Hierarchical organization of speech perception in human auditory cortex.
635 Frontiers in neuroscience, 8:406, 2014.
- 636 [32] Eric J Hunter, Jan G Švec, and Ingo R Titze. Comparison of the produced
637 and perceived voice range profiles in untrained and trained classical singers.
638 Journal of Voice, 20(4):513–526, 2006.
- 639 [33] Jintao Jiang, Marcia Chen, and Abeer Alwan. On the perception of voicing
640 in syllable-initial plosives in noise. The Journal of the Acoustical Society
641 of America, 119(2):1092–1105, 2006.
- 642 [34] Eric R Kandel, James H Schwartz, Thomas M Jessell, Steven Siegelbaum,
643 A James Hudspeth, and Sarah Mack. Principles of neural science, volume 4.
644 McGraw-hill New York, 2000.
- 645 [35] John Kominek, Tanja Schultz, and Alan W Black. Synthesizer voice quality
646 of new languages calibrated with mean mel cepstral distortion. In Spoken
647 Languages Technologies for Under-Resourced Languages, 2008.
- 648 [36] Sergey Korolev, Amir Safiullin, Mikhail Belyaev, and Yulia Dodonova.
649 Residual and plain convolutional neural networks for 3d brain mri classifi-

- 650 cation. In 2017 IEEE 14th international symposium on biomedical imaging
651 (ISBI 2017), pages 835–838. IEEE, 2017.
- 652 [37] Joseph G Makin, David A Moses, and Edward F Chang. Machine transla-
653 tion of cortical activity to text with an encoder–decoder framework. Nature
654 Neuroscience, 23(4):575–582, 2020.
- 655 [38] Brian T Miller and Mark D’Esposito. Searching for “the top” in top-down
656 control. Neuron, 48(4):535–538, 2005.
- 657 [39] Muge Ozker, Werner Doyle, Orrin Devinsky, and Adeen Flinker. Cortical
658 network underlying speech production during delayed auditory feedback.
659 bioRxiv, 2021.
- 660 [40] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna
661 Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations
662 from deep networks via gradient-based localization. In Proceedings of the
663 IEEE international conference on computer vision, pages 618–626, 2017.
- 664 [41] Jennifer Shum, Lora Fanda, Patricia Dugan, Werner K Doyle, Orrin Devin-
665 sky, and Adeen Flinker. Neural correlates of sign language production re-
666 vealed by electrocorticography. Neurology, 95(21):e2880–e2889, 2020.
- 667 [42] Kristina Simonyan, Hermann Ackermann, Edward F Chang, and Jeremy D
668 Greenlee. New developments in understanding the complexity of human
669 speech production. Journal of Neuroscience, 36(45):11440–11448, 2016.
- 670 [43] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin
671 Wattenberg. Smoothgrad: removing noise by adding noise. arXiv preprint
672 arXiv:1706.03825, 2017.
- 673 [44] Donald T Stuss and Robert T Knight. Principles of frontal lobe function.
674 Oxford University Press, 2013.

- 675 [45] Cees H Taal, Richard C Hendriks, Richard Heusdens, and Jesper Jensen.
676 A short-time objective intelligibility measure for time-frequency weighted
677 noisy speech. In 2010 IEEE international conference on acoustics, speech
678 and signal processing, pages 4214–4217. IEEE, 2010.
- 679 [46] Ran Wang, Xupeng Chen, Amirhossein Khalilian-Gourtani, Zhaoxi Chen,
680 Leyao Yu, Adeen Flinker, and Yao Wang. Stimulus speech decoding from
681 human cortex with generative adversarial network transfer learning. In
682 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI),
683 pages 390–394. IEEE, 2020.
- 684 [47] Ran Wang, Yao Wang, and Adeen Flinker. Reconstructing speech stim-
685 ulti from human auditory cortex activity using a WaveNet approach. In
686 2018 IEEE Signal Processing in Medicine and Biology Symposium (SPMB),
687 pages 1–6. IEEE, 2018.
- 688 [48] Andrew I Yang, Xiuyuan Wang, Werner K Doyle, Eric Halgren, Chad
689 Carlson, Thomas L Belcher, Sydney S Cash, Orrin Devinsky, and Thomas
690 Thesen. Localization of dense intracranial electrode arrays using magnetic
691 resonance imaging. Neuroimage, 63(1):157–165, 2012.
- 692 [49] James F.A. Poulet and Berthold Hedwig. The cellular basis of a corollary
693 discharge. Science, 311: 518–522, 2006.
- 694 [50] Eliades, Steven J., and Xiaoqin Wang. Neural substrates of vocalization
695 feedback monitoring in primate auditory cortex. Nature, 453(7198):1102-
696 1106, 2008.

697 **5 Acknowledgements**

698 We would like to thank Robert Knight and Sasha Devore for providing helpful
699 comments on the manuscript. This work was supported by the National Sci-
700 ence Foundation under Grant No. IIS-1912286 (Y.W. and A.F.) and National
701 Institute of Health R01NS109367 (to A.F.).

702 **6 Author contributions**

703 R.W. conceived and implemented the decoding algorithm and interpreted the
704 model with assistance from Y.W. and A.F.; X.C participated in the data pro-
705 cessing and performance evaluation; A.K.G participated in the data processing
706 and visualization; L.Y. participated in data acquisition, preprocessing, and vi-
707 sualization; P.D. and D.F. provided clinical care; W.D. provided neurosurgical
708 clinical care; O.D. assisted in patient care and consent; Y.W. led the research
709 project and advised from engineering perspective; A.F. co-led the project with
710 Y.W., participated in all data acquisition, and advised from neuroscience per-
711 spective; R.W. and A.F. co-wrote the manuscript with input from all authors.

712 **7 Competing interests**

713 The authors declare no competing financial interests.

714 A Extended Data

715 A.1 Additional Decoding Framework Details

716 A.1.1 Differentiable speech synthesizer

717 In a traditional vocoder, speech is generated by switching between voiced and
718 unvoiced content. Each content comes from an autoregressive system driven by
719 a certain excitation signal that is either a harmonic signal or a white noise sig-
720 nal [11]. Inspired by such a process, we construct our speech synthesizer shown
721 in Fig. 6. It consists of two pathways. The *voice pathway* generates a voiced
722 component by driving a harmonic excitation with time-varying fundamental fre-
723 quency (i.e., pitch) f_0^t through a voice filter consisting of N formant filters, each
724 described by a center frequency f_i^t and an amplitude $a_i^t, i = 1, 2, \dots, N$. Note
725 that we parameterize the bandwidth b_i^t as a function of the center frequency f_i^t ,
726 as discussed later. The *unvoice pathway* generates an unvoiced component by
727 driving a white noise through an unvoice filter described as a center frequency
728 f_n^t , bandwidth b_n^t and amplitude a_n^t (in addition to the N formant filters for
729 the voice pathway). These two components are adaptively combined with a
730 time-varying mixing factor α^t , controlling the relative contribution between
731 voiced sounds (for sonorant phonemes including vowels and nasals) and un-
732 voiced sounds (for voiceless plosives and fricatives such as /p/, /s/). The voiced
733 plosives and fricatives (such as /b/, /z/) can be generated as a combination of
734 voiced and unvoiced components. Finally, the combined signal is amplified by a
735 loudness parameter L_t . In our study, we used $N = 6$ formants. The synthesizer
736 is driven by a total of 18 time-varying speech parameters, including the fun-
737 damental (or pitch) frequency f_0^t , the mixing factor between the two pathways
738 α^t , the 12 parameters for the voice filter (f_i^t, a_i^t) and the three parameters for
739 the unvoice filter f_n^t, b_n^t, a_n^t , and the loudness L^t . Given the parameter values

740 at each time sample, the synthesizer can generate a spectrogram sample. The
741 spectrogram is a differentiable function of the speech parameters so that we can
742 back-propagate the gradient of the training loss in terms of the predicted spec-
743 trogram to the speech parameters, which can then be backpropagated to either
744 the speech encoder or the ECoG decoder parameters. Specifically, let the $V^t(f)$
745 represent the spectrogram of the voicing component, $U^t(f)$ that of the unvoicing
746 component, and $\alpha^t \in [0, 1]$ the mixing factor. The combined spectrogram can
747 be written as $S^t(f) = \alpha^t V^t(f) + (1 - \alpha^t) U^t(f)$. Finally, the synthesized speech
748 spectrogram is $\tilde{S}^t(f) = L^t S^t(f)$, where L^t is the loudness that modulates the
749 signal cross time.

750 **Formant filters in the voice pathway** The filter in the voice pathway
751 consists of multiple formant filters, corresponding to the multiple formants as-
752 sociated with vowels. The formant filter shape over frequency, which is related
753 to the resonance property of the vocal tract, is closely related to the timbre
754 of speakers' voice [32]. We have found that a predefined analytic form such as
755 generalized Gaussian cannot cover all feasible filter shapes. Instead, we learn
756 a speaker-dependent prototype filter for each formant based on the speaker's
757 natural speech. We represent the prototype filter ($G_i(f)$) for the i -th formant
758 as a piecewise linear function, linearly interpolated from $g^i[m], m = 1 \dots M$, the
759 amplitudes of the filter at M uniformly sampled frequencies up to f_{max} . We
760 restrict the resulting filter $G_i(f)$ to be unimodal (with a single peak of value
761 1) by properly constraining $g[m]$. Given $g[m], m = 1 \dots M$, the peak frequency
762 f_i^{proto} and the half-power bandwidth b_i^{proto} can be determined. The actual for-
763 mant filter at any time can be written as a shifted and scaled version of $G_i(f)$.
764 Specifically, at time t , given an amplitude (a_i^t), a center frequency (f_i^t), and a

765 bandwidth (b_i^t), the i -th formant filter is given by

$$F_i^t(f) = a_i^t \cdot G_i \left(\frac{b_i^{proto}}{b_i^t} \cdot (f - f_i^t) + f_i^{proto} \right) \quad (1)$$

766 Then the filter for the voice pathway with N formant filters can be written
767 as $F_h^t(f) = \sum_{i=1}^N F_i^t(f)$. We learn the parameters $g[m], m = 1 \dots M$ for $G_i(f)$
768 during the unsupervised pre-training of the speech encoder, which does not
769 require neural data. Fitting such a prototype filter is not data-hungry even
770 with a relatively large M . We used $M = 20$ in our experiment. Although two
771 formants ($N=2$) have been shown to suffice for intelligible reconstruction [7],
772 we use $N=6$ in our experiments for more accurate synthesis. We denote the
773 parameter set for the voice filter at time t by $\mathcal{S}^t = \{(f_i^t, a_i^t, b_i^t) | i \in \{1, \dots, N\}\}$.
774 As explained later, the bandwidth b_i^t parameters are not independent speech
775 parameters, rather functions of the center frequencies f_i^t .

776 **Unvoice filter** For the unvoice pathway, we add a broadband filter described
777 by $\{(f_n^t, a_n^t, b_n^t)\}$. The shape of this filter $F_n^t(f)$ follows equation (1) but with
778 the filter coefficients $(\alpha_i^t, f_i^t, b_i^t)$ replaced by $(\alpha_n^t, f_n^t, b_n^t)$. The bandwidth is
779 constrained to satisfy $b_n^t > 2000\text{Hz}$, following the broadband nature of obstruent
780 phonemes. We also keep the multiple formant filters in the voice filter described
781 by \mathcal{S}^t . This is motivated by the fact that human beings differentiate consonants
782 with similar sounds such as /p/ and /d/, not only by the immediate burst
783 of these sounds, but also the development of the following formant frequency
784 until the next vowel [33]. To encode such formant transitions, we use the same
785 formant filter parameters for modeling the narrow passbands in both the voiced
786 component and the unvoiced component. The parameter set for the unvoiced
787 component is thus $T^t = \mathcal{S}^t \cup \{(f_n^t, a_n^t, b_n^t)\}$. The overall filter for the unvoice
788 pathway is: $F_n^t(f) = F_n^t(f) + \sum_{i=1}^N F_i^t(f)$.

789 To further reduce the parameter space dimension, we model the bandwidth
790 b_i^t of a formant filter as a piecewise linear function of the center frequency f_i^t .

791 We assume

$$b_i^t = \begin{cases} a(f_i^t - f_\theta) + b_0, & \text{if } f_i^t > f_\theta \\ b_0, & \text{otherwise} \end{cases}$$

792 where threshold frequency f_θ , slope a , and baseline bandwidth b_0 are three
793 parameters that can be learnt during unsupervised pre-training, shared among
794 all formant filters.

795 **Harmonic excitation** In the voice pathway, the voice filter is applied on
796 the harmonic excitation. This pathway models the human production of vowels
797 and nasals, which results from the voice excited by the vocal cord shaped by the
798 vocal tract. The excitation is constructed by sinusoidal harmonic oscillations
799 with a time varying fundamental frequency f_0^t . Inspired by the formulation
800 in [12], we define the harmonic excitation h^t as: $h^t = \sum_{k=1}^K h_k^t$, where K is the
801 total number of harmonics ($K=80$ in our experiment). Assuming the initial
802 phase is 0, each harmonic resonance h_k^t at time step t has an instant phase that
803 is the accumulation of resonance frequency in the past. Specifically, the k -th
804 resonance at time step t is $h_k^t = \sin(2\pi \sum_{\tau=0}^t f_k^{(\tau)})$, where $f_k^{(t)} = k f_0^{(t)}$. Denoting
805 the spectrogram of h^t as $H^t(f)$, the spectrogram of the voice component is the
806 multiplication of $H^t(f)$ and the voice filter, i.e., $V^t(f) = H^t(f) \circ F_h^t(f)$.

807 **Noise excitation** The unvoice pathway models consonants like plosives and
808 fricatives, where the vocal tract and human mouth filter the airflow through
809 the mouth. It follows a similar process as in the harmonic counterpart. The
810 major difference is that the excitation being filtered becomes stationary white
811 Gaussian distributed noise $\hat{n}(t) \sim \mathcal{N}(0, 1)$, with a corresponding spectrogram
812 $N^t(f)$. The filtered noise spectrogram (i.e., the unvoice component) is $U^t(f) =$

813 $N^t(f) \circ F_n^t(f)$.

814 **A.1.2 ECoG decoder and speech encoder**

815 The ECoG decoder is constructed by a three-dimensional ResNet that treats
816 time-varying signals on an ECoG grid array as spatiotemporal three-dimensional
817 tensors (width \times height \times time duration). As is depicted in Figure 6c, after an
818 initial temporal convolutional layer (with 128 feature map filters and a kernel
819 size of $1 \times 1 \times 9(72ms)$), the signal passes through eight residual blocks. Each
820 block contains two three-dimensional convolutional layers (with 128 feature map
821 filters, each has kernel size of $3 \times 3 \times 5(40ms)$). The output of the residual blocks
822 creates a shared latent representation consisting of 128 feature maps (each is a
823 one-dimensional temporal signal by average pooling the two spatial dimensions),
824 which is then fed into different output heads (each applies each consists of one or
825 two fully connected layers acting on the 128 features at the same time point) to
826 generate speech parameters. The overall temporal receptive field for generating
827 one speech parameter sample is 73 temporal samples of 584 ms.

828 The speech encoder network architecture we choose is as simple as possible
829 to demonstrate the effectiveness of the speech synthesizer design. In the exper-
830 iment, we use three layers of temporal convolution (we treat the frequency axis
831 of the spectrogram as the feature dimension) to generate a latent representa-
832 tion (Figure 6d). Each convolutional layer has 128 feature maps and a temporal
833 kernel size of 3 frames (24ms). To output the speech parameter, we apply the
834 same multi-head structure to the latent representation as in the last layer of the
835 ECoG decoder.

836 **A.1.3 Loss and training hyper-parameters**

The speech encoder is trained with a weighted average of the mixed spectral loss and the parameter loss. The mixed spectral loss [12] is defined as:

$$L_{MSS}(\tilde{S}^t(f), S^t(f)) = L_{\text{lin}}(\tilde{S}^t(f), S^t(f)) + L_{\text{mel}}(\tilde{S}^t(f), S^t(f)),$$

837 in which,

$$\begin{aligned} L_{\text{lin}}(x, y) &= \|x - y\|_1 + \|\log x - \log y\|_1 \\ L_{\text{mel}}(x, y) &= \|x_{\text{mel}} - y_{\text{mel}}\|_1 + \|\log x_{\text{mel}} - \log y_{\text{mel}}\|_1 \end{aligned}$$

838 where $S^t(f)$ and $\tilde{S}^t(f)$ denote the ground truth and reconstructed spectrograms,
839 respectively, subscript *lin* means that the frequency is in the linear scale while
840 the subscript *mel* means the frequency is in the mel scale. In our experiments,
841 we use 256 frequency samples (ranging from 0-8000 Hz) for both linear scale
842 and mel scale speech spectrograms.

Let's denote the j -th reconstructed speech parameter as \tilde{P}_j^t and its reference P_j^t , the overall training loss for the ECoG decoder becomes:

$$\begin{aligned} L &= L_{\text{spectrogram}} + L_{\text{speechparameters}} \\ &= \lambda_0 L_{MSS}(\tilde{S}^t(f), S^t(f)) + \sum_j \lambda_j \left(\left\| \tilde{P}_j^t - P_j^t \right\|_2^2 \right) \end{aligned}$$

843 where λ_j balance the contribution from different loss terms since they have
844 different physical meanings and scales.

845 Both the speech encoder and ECoG decoder are fitted by Adam optimizer
846 with hyper-parameters: $lr = 10^{-3}, \beta_1 = 0.9, \beta_2 = 0.999$. We train an individual
847 ECoG decoder and speech encoder per patient. The pre-training of the speech

848 encoder and the training of the ECoG decoder share the same training/testing
849 set partition.

850 B Additional Figures and Tables

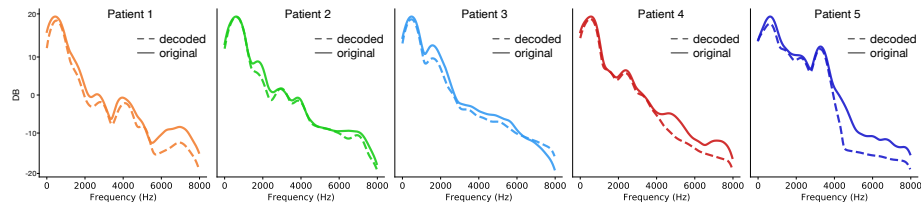


Figure E1: The spectral energy distribution of the decoded and original speech for five patients. Visualized by averaging the broadband spectrograms magnitude across time of all test samples.

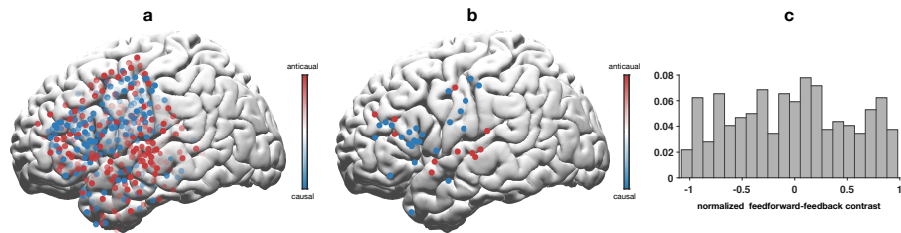


Figure E2: Normalized contrast of feedforward vs. feedback contribution. (a) electrode level feedforward-feedback contribution contrast, normalized by the sum of feedforward and feedback contribution magnitude. (b) electrodes with large feedforward-feedback polarity with the normalized contrast magnitude > 0.9. (c) The histogram of the normalized contrast. Positive bins correspond to anticausal polarization.

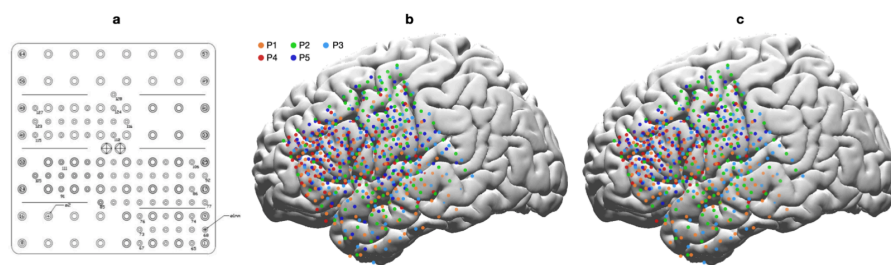


Figure E3: Electrodes array and implant location of all five patients (P1-P5) in our experiments. (a) The 128 electrodes on the hybrid density ECoG array. (b) All electrodes on cortex (MNI). (c) All electrodes with usable data. Only data from these electrodes are used to train the ECoG decoder models.

Anatomical region	p	z
rSTG	0.0332	-2.9628
cSTG	1.607E-15	9.6234
rMTG	2.5108E-04	4.9359
mMTG	1.5257E-13	9.0185
cMTG	0.2269	1.5656
ventralprecentral	4.9511E-8	-7.1409
dorsalprecentral	0.4349	0.6525
postcentral	6.419E-04	-4.9612
cMFG	0.0248	3.1417
rMFG	0.1988	1.7202
parstriangularis	2.6715E-06	6.3518
parsopercularis	8.0693E-15	-9.6185
supramarginal	1.1144E-04	5.3919

Table 1: Statistics of data in Figure 3f. The P-value and Z-value are reported for Wilcoxon sign rank test between feedback and feedforward contributions across all electrodes and test trials within each anatomical region. The Z-value represents the rank based test statistic with positive values reflecting anticausal contributions and negative values reflecting causal contributions.

Anatomical region	Causal pre		Causal during		Anticausal	
	p-value	z-value	p-value	z-value	p-value	z-value
cSTG	0.7226	0.075	5.6745E-10	8.4078	1.2168E-22	10.9874
rSTG	0.4942	-0.0510	1.2703E-07	6.8284	0.1557	1.9149
mMTG	5.3721E-06	6.2146	0.3689	0.2741	1.8216E-28	12.1658
cMTG	0.1671	1.0126	0.231	-0.501	0.4781	-0.3012
rMTG	5.1253E-19	10.1547	0.1293	2.1082	1.4923E-10	8.2051
ventralprecentral	1.7845E-58	16.2047	3.0286E-77	17.5451	2.2394E-60	17.1839
dorsalprecentral	2.9083E-12	8.0932	8.9452E-04	4.4590	1.4512E-09	7.9235
postcentral	3.67853E-91	21.4986	9.34051E-104	22.1393	6.9834E-34	14.0134
supramarginal	2.2905E-06	6.7810	0.5924	-0.2945	1.8542E-07	6.9384
paropercularis	3.9368E-76	19.0572	3.843E-72	18.5329	3.083E-04	5.3823
parstriangularis	7.2744E-77	19.5782	5.8573E-31	13.9374	2.0273E-37	14.4676
rMFG	2.3846E-27	12.2940	2.0371E-07	7.8460	0.3643	0.3823
cMFG	4.0274E-26	11.0042	2.83632E-07	6.9027	9.02834E-19	9.1881

Table 2: Statistics of data in Figure 4 and 5. Per anatomical region P-value and Z-value are reported for Wilcoxon sign rank test between the each regions' contribution and the shuffled model's contribution (control curves in Figure 5). The Z-value represents the rank based test statistic with positive values reflecting larger real contributions compared with shuffled contributions. This is shown for the causal model (pre-production period), causal model (during production period), and anticausal model, respectively. Curves of each individual electrode and test trial are considered one sample, and are averaged across time to perform the Wilcoxon sign rank test. The red marked regions in the table are highlighted to show no significance (P-value>0.05) and are omitted when plotting the curves in Figure 5 as described (Method - Revealing delay-dependent contribution of different cortical regions from the trained ECoG to speech model - Visualizing per region temporal contribution receptive field).

Anatomical region	Causal vs. Anticausal		Causal during vs pre	
	P-value	Z-value	P-value	Z-value
cSTG	2.6789E-17	9.6711	4.718E-04	3.696
rSTG	0.0343	-2.9457	6.2075E-04	4.7427
mMTG	3.2252E-13	9.0928	4.5863E-04	-4.0475
cMTG	0.3930	1.0021	0.2718	-1.1957
rMTG	1.8511E-04	5.1625	1.0173E-10	-8.9283
ventralprecentral	2.8012E-15	-10.0562	8.2757E-05	5.0475
dorsalprecentral	0.6492	0.2967	5.5615E-04	-3.4394
postcentral	3.0581E-08	-6.1286	0.3037	1.7462
supramarginal	1.9928E-07	6.0301	4.8257E-06	-6.0274
parsopercularis	8.6228E-18	-10.0274	0.5922	0.1582
parstriangularis	0.0162	3.9003	3.2532E-32	-12.4583
rMFG	0.0021	-4.9475	2.5714E-04	-5.0131
cMFG	0.0045	3.9862	3.0747E-09	-7.0652

Table 3: Statistics of data in Figure 4 and 5. Per anatomical region P-value and Z-value are reported for Wilcoxon sign rank test between the causal (during production period) model and the anticausal model (The positive/negative Z-values represent the direction of the contribution where positive values denote anticausal greater than causal), as well as the causal model between during- and pre- epochs (The positive/negative Z-values represent the direction of the contribution where positive values denote during production greater than pre-production). Curves of each individual electrode and test trial are considered as one sample, and are averaged across the time epoch to perform the Wilcoxon sign rank test. The red marked regions in the table are highlighted to denote no significance (P-value>0.05).

Anatomical region	Causal (pre)	Causal (during)	Anticausal
cSTG	-	-352	168
rSTG	-	-256	-
mMTG	-176	-	240
cMTG	-	-	-
rMTG	-192	-	312
ventralprecentral	-196	-208	280
dorsalprecentral	-192	-184	144
postcentral	-248	-256	192
supramarginal	-120	-	184
parsopercularis	-248	-280	232
parstriangularis	-240	-336	264
rMFG	-248	-304	-
cMFG	-248	-304	208

Table 4: Peak time of each anatomical region curves in Figure 5 a,b,c. Each column reports the peak time of the temporal receptive field curves for the causal model (pre-production), causal model (during production), and anticausal model, respectively. The peak of each region is calculated based on the averaged curve (averaged across trials and electrodes within the region).