

1 **Towards reconstructing intelligible speech from the human auditory**
2 **cortex**

3
4
5
6

7 Hassan Akbari^{1,2}, Bahar Khalighinejad^{1,2}, Jose L. Herrero^{3,4}, Ashesh D. Mehta^{3,4}, Nima Mesgarani^{1,2}

8 ¹*Mortimer B. Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY*

9 ²*Department of Electrical Engineering, Columbia University, New York, NY*

10 ³*Hofstra Northwell School of Medicine, Manhasset, NY, United States*

11 ⁴*The Feinstein Institute for Medical Research, Manhasset, NY, United States*

12
13
14
15

16 Correspondence to: nima@ee.columbia.edu

17

18 **Abstract**

19 Auditory stimulus reconstruction is a technique that finds the best approximation of the acoustic stimulus
20 from the population of evoked neural activity. Reconstructing speech from the human auditory cortex
21 creates the possibility of a speech neuroprosthetic to establish a direct communication with the brain and
22 has been shown to be possible in both overt and covert conditions. However, the low quality of the
23 reconstructed speech has severely limited the utility of this method for brain-computer interface (BCI)
24 applications. To advance the state-of-the-art in speech neuroprosthesis, we combined the recent
25 advances in deep learning with the latest innovations in speech synthesis technologies to reconstruct
26 closed-set intelligible speech from the human auditory cortex. We investigated the dependence of
27 reconstruction accuracy on linear and nonlinear (deep neural network) regression methods and the
28 acoustic representation that is used as the target of reconstruction, including auditory spectrogram and
29 speech synthesis parameters. In addition, we compared the reconstruction accuracy from low and high
30 neural frequency ranges. Our results show that a deep neural network model that directly estimates the
31 parameters of a speech synthesizer from all neural frequencies achieves the highest subjective and
32 objective scores on a digit recognition task, improving the intelligibility by 65% over the baseline method
33 which used linear regression to reconstruct the auditory spectrogram. These results demonstrate the
34 efficacy of deep learning and speech synthesis algorithms for designing the next generation of speech BCI
35 systems, which not only can restore communications for paralyzed patients but also have the potential to
36 transform human-computer interaction technologies.

37

38 **Introduction**

39 Auditory stimulus reconstruction is an inverse mapping technique that finds the best approximation of
40 the acoustic stimulus from the population of evoked neural activity. Stimulus reconstruction was originally
41 proposed as a method to study the representational properties of the neural population¹⁻⁵ because this
42 method enables the intuitive interpretation of the neural responses in the stimulus domain.
43 Reconstructing speech from the neural responses recorded from the human auditory cortex⁶, however,
44 opens up the possibility of using this technique as a speech brain-computer interface (BCI) to restore
45 speech in severely paralyzed patients (for a review, see these references⁷⁻⁹). The ultimate goal of a speech
46 neuroprosthesis is to create a direct communication pathway to the brain with the potential to benefit
47 patients who have lost their ability to speak, which can result from a variety of clinical disorders leading
48 to conditions such as locked-in syndrome^{10,11}. The practicality of using speech decoding methods in a
49 neuroprosthetic device to restore speech communication was further supported by studies showing

50 successful decoding of speech during both overt and covert (imagined) conditions^{12–16}. These studies
51 showed successful decoding of imagined articulations^{13,14}, imagined word repetition¹⁵, and silent reading
52 of speech¹⁶ from auditory cortical areas, including the superior temporal gyrus (STG). While previous
53 studies have established the feasibility of reconstructing speech from neural data, the quality of the
54 reconstructed audio so far has been too low to merit subjective evaluation. For this reason, the
55 reconstructed sounds in previous studies have been evaluated only using objective measures such as
56 correlation or recognition accuracy^{3,6,8,13,17–25}. The low quality of the reconstructed sound is currently a
57 major limiting factor in actualizing speech BCI systems⁷.

58 The acoustic representation of the stimulus that is used as the decoding target can significantly
59 impact the quality and accuracy of reconstructed sounds. Previous studies have used magnitude
60 spectrogram (time-frequency representation)^{3,20}, speech envelope^{21,22}, spectrotemporal modulation
61 frequencies^{6,13,23}, and discrete units such as phonemes and phonetic categories^{8,17,24,25} and words^{18,19}.
62 Using discrete units can be advantageous by allowing for discriminative training. However, decoding
63 discrete representations of speech such as phonemes eliminates the paralinguistic information such as
64 speaker features, emotion, and intonation. In comparison, reconstructing continuous speech provides the
65 possibility of real-time, continuous feedback that can be delivered to the user to promote coadaptation
66 of the subject and the BCI algorithm^{26,27} for enhanced accuracy. A natural choice is to directly estimate
67 the parameters of a speech synthesizer from neural data, but this has not been attempted previously
68 because the process requires a highly accurate estimation of several vocoder parameters, which is hard
69 to achieve with traditional machine-learning techniques.

70 To advance the state-of-the-art in speech neuroprosthesis, we aimed to increase the intelligibility
71 of the reconstructed speech by combining recent advances in deep learning²⁸ with the latest innovations
72 in speech synthesis technologies. Deep learning models have recently become the dominant technique
73 for acoustic and audio signal processing^{29–32}. These models can improve reconstruction accuracy by
74 imposing more complete constraints on the reconstructed audio by better modeling the statistical
75 properties of the speech signal³. At the same time, nonlinear regression can invert the nonlinearly
76 encoded speech features in neural data^{33,34} more accurately.

77 We examined the effect of three factors on the reconstruction accuracy: 1) the regression
78 technique (linear regression versus nonlinear deep neural network), 2) the representation of the speech
79 intended for reconstruction (auditory spectrogram versus speech vocoder parameters), and 3) the neural
80 frequency range used for regression (low frequency versus high-gamma envelope) (Fig. 1A). Our results
81 showed that a deep neural network model that uses all neural frequencies to directly estimate the

82 parameters of a speech vocoder achieves the highest subjective and objective scores, both for
83 intelligibility and the quality of reconstruction in a digit recognition task. These results represent an
84 important step toward successful implementation of the next generation of speech BCI systems.

85

86 **Results**

87 **Neural recordings:** We used invasive electrocorticography (ECoG) to measure neural activity from five
88 neurosurgical patients undergoing treatment for epilepsy as they listened to continuous speech sounds.
89 Two of the five subjects had high-density subdural grid electrodes implanted in the left hemisphere with
90 coverage primarily over the superior temporal gyrus (STG), and four of the five subjects had depth
91 electrodes with coverage of Heschl's gyrus (HG). All subjects had self-reported normal hearing. Subjects
92 were presented with short continuous stories spoken by four speakers (two females, total duration: 30
93 minutes). To ensure that the subjects were engaged in the task, the stories were randomly paused, and
94 the subjects were asked to repeat the last sentence.

95 The test data consisted of continuous speech sentences and isolated digit sounds. We used eight
96 sentences (40 seconds total) to evaluate the objective quality of the reconstruction models. The sentences
97 were repeated six times in random order, and the neural data was averaged over the six repetitions to
98 reduce the effect of neural noise on comparison of reconstruction models (see Supp. Fig. 1 for the effect
99 of averaging). The digit sounds were used for subjective intelligibility and quality assessment of
100 reconstruction methods and were taken from a publicly available corpus, TI-46³⁵. We chose 40 digit
101 sounds (zero to nine), spoken by four speakers (two females) that were not included in the training of the
102 models. Reconstructed digits were used as the test set to evaluate subjective intelligibility and quality of
103 the models. Two ranges of neural frequencies were used in the study. Low-frequency (0–50 Hz)
104 components of the neural data were extracted by filtering the neural signals using a lowpass filter. The
105 high-gamma envelope³⁶ was extracted by filtering the neural signals (70 to 150 Hz) and calculating the
106 Hilbert envelope³⁷.

107

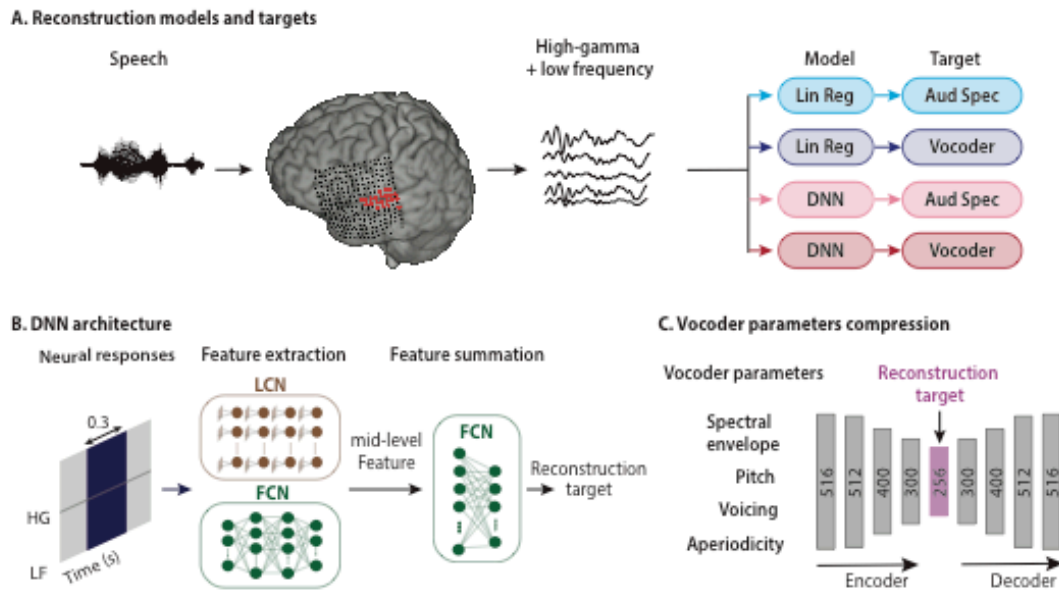
108

109

110 **Regression models:** The input to the regression models was a sliding window over the neural data with a
111 duration of 300 ms (Fig. 1B), and the hop size of 10 ms. The duration of the sliding window was chosen to
112 maximize reconstruction accuracy (supp. Fig. 2). We compared the performance of linear and nonlinear
113 regression models to reconstruct the stimulus from the neural signals. The linear regression finds a linear
114 mapping between the response of a population of neurons to the stimulus representation^{3,6}. This method
115 effectively assigns a spatiotemporal filter to each electrode estimated by minimizing the mean-squared-
116 error (MSE) between the original and reconstructed stimulus.

117 The nonlinear regression model was implemented using a deep neural network (DNN). We
118 designed a deep neural network architecture with two stages: 1) feature extraction and 2) feature
119 summation networks³⁸⁻⁴⁰ (Fig. 1B). In this framework, a high-dimensional representation of the input
120 (neural responses) is first calculated, which results in mid-level features (output of the feature extraction
121 network). These mid-level features are then input to the feature summation network to regress the
122 output of the model (acoustic representation). The feature summation network in all cases was a two-
123 layer fully connected network (FCN) with regularization, dropout⁴¹, batch normalization⁴², and
124 nonlinearity between each layer. For feature extraction, we compared the efficacy of five different
125 network architectures for auditory spectrogram and vocoder reconstruction (Methods, Supp. Table 1 for
126 details of each network). Specifically, we found that the fully connected network (FCN), in which no
127 constraint was imposed on the connectivity of the nodes in each layer of the network to the previous
128 layer, achieved the best performance for reconstructing the auditory spectrogram. However, the
129 combination of the FCN and a locally connected network (LCN), which constrains the connectivity of each
130 node to only a subset of nodes in the previous layer, achieved the highest performance for the vocoder
131 representation (Supp. Tables 4, 5). In the combined FCN+LCN, the outputs of the two parallel networks
132 are concatenated and used as the mid-level features (Fig.1B).

133
134
135
136
137
138
139
140



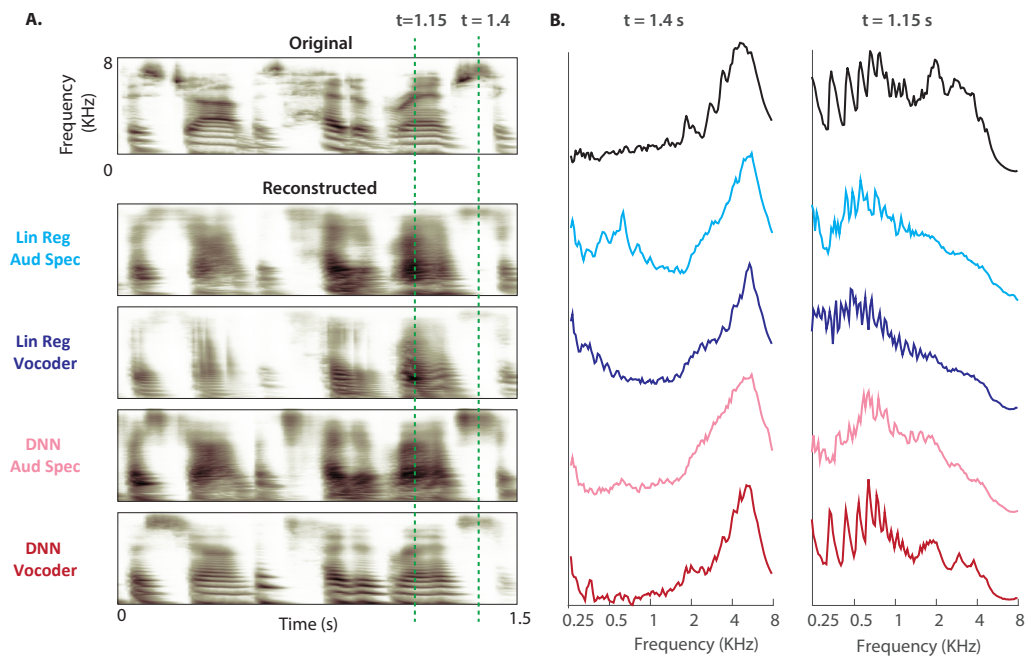
141
 142 **Figure 1. Schematic of the speech reconstruction method.** (A) Subjects listened to natural speech
 143 sentences. The population of evoked neural activity in the auditory cortex of the listener was then used
 144 to reconstruct the speech stimulus. The responsive electrodes in an example subject are shown in red.
 145 High and low frequency bands were extracted from the neural data. Two types of regression models and
 146 two types of speech representations were used, resulting in four combinations: linear regression to
 147 auditory spectrogram (light blue), linear regression to vocoder (dark blue), DNN to auditory spectrogram,
 148 and DNN to vocoder (dark red). (B) The input to all models was a 300 ms sliding window containing both
 149 low frequency (LF) and the high-gamma envelope (HG). The DNN architecture consists of two modules:
 150 feature extraction and feature summation networks. Feature extraction for auditory spectrogram
 151 reconstruction was a fully connected neural network (FCN). For vocoder reconstruction, the feature
 152 extraction network consisted of an FCN concatenated with a locally connected network (LCN). The feature
 153 summation network is a two-layer fully connected neural network (FCN). (C) Vocoder parameters consist
 154 of spectral envelope, fundamental frequency (f_0), voicing, and aperiodicity (total of 516 parameters). An
 155 autoencoder with a bottleneck layer was used to reduce the 516 vocoder parameters to 256. The
 156 bottleneck features were then used as the target of reconstruction algorithms. The vocoder parameters
 157 were calculated from the reconstructed bottleneck features using the decoder part of the autoencoder
 158 network.
 159
 160

161 **Acoustic representations:** We used two types of acoustic representation of the audio as the target for
162 reconstruction: auditory spectrogram and speech vocoder. The auditory spectrogram was calculated
163 using a model of the peripheral auditory system^{43,44}, which estimates a time-frequency representation of
164 the acoustic signal on a tonotopic frequency axis. The reconstruction of the waveform from the auditory
165 spectrogram is achieved using an iterative convex optimization procedure⁴³ because the phase of the
166 signal is lost during this procedure.

167 For speech vocoder, we used a vocoder-based, high-quality speech synthesis algorithm
168 (WORLD⁴⁵), which synthesizes speech from four main parameters: 1) spectral envelope, 2) f_0 or
169 fundamental frequency, 3) band aperiodicity, and 4) a voiced-unvoiced (VUV) excitation label (Fig. 1C).
170 These parameters are then used to re-synthesize the speech waveform. This model can reconstruct high-
171 quality speech and has been shown to outperform other methods including STRAIGHT⁴⁶. The large
172 numbers of the parameters in the vocoder (516 total) and the susceptibility of the synthesis quality on
173 inaccurate estimation of parameters however pose a challenge. To remedy this, we first projected the
174 sparse vocoder parameters onto a dense subspace in which the number of parameters can be reduced,
175 which allows better training with a limited amount of data. We used a dimensionality reduction technique
176 that relies on an autoencoder (AEC)⁴⁷(Fig. 1C), which compresses the vocoder parameters into a smaller
177 space (encoder, 256 dimensions, Supp. Table 3) and subsequently recovers (decoder) the original vocoder
178 parameters from the compressed features (Fig. 1C). The compressed features (also called bottleneck
179 features) are used as the target for the reconstruction network. By adding noise to the bottleneck features
180 before feeding them to the decoder during training, we can make the decoder more robust to unwanted
181 variations in amplitude, which is necessary due to the noise inherently present in the neural signals. The
182 autoencoder was trained on 80 hours of speech using a separate speech corpus (Wall Street Journal I⁴⁸).
183 During the test phase, we first reconstructed the bottleneck features from the neural data, and
184 subsequently estimated the vocoder parameters using the decoder part of the autoencoder (Fig. 1C). The
185 reconstruction accuracy of individual vocoder parameters with a neural network shows varied
186 improvement over the linear model, where pitch estimation is improved the most (%157.2), followed by
187 aperiodicity (%18.5), spectral envelope (%6.2), and voiced-unvoiced parameter (%0.15, Supp. Fig. 3).

188 Figure 2B shows the example reconstructed auditory spectrograms from each of the four
189 combinations of the regression models (linear regression and DNN) and acoustic representation (auditory
190 spectrogram and vocoder). Comparison of the auditory spectrograms in Figure 2A shows that 1) the
191 overall frequency profile of the speech utterance is better preserved by the DNN compared to the linear
192 regression model, and 2) the harmonic structure of speech is recovered only in the DNN-Vocoder model.

193 These observations are shown more explicitly in Figure 2B, where the magnitude power of frequency
194 bands is shown during an unvoiced ($t = 1.4$ sec) and a voiced speech sound ($t = 1.15$ sec, shown with
195 dashed lines). The frequency profile of original and reconstructed auditory spectrograms during the
196 unvoiced sound shows a more accurate reconstruction of low and high frequencies for the DNN models
197 (Fig. 2B left, comparison of blue and red plots). The comparison of frequency profiles during the voiced
198 sound (Fig. 2B, right) reveals the recovery of speech harmonics only in the DNN-Vocoder model
199 (comparison of top and bottom plots).



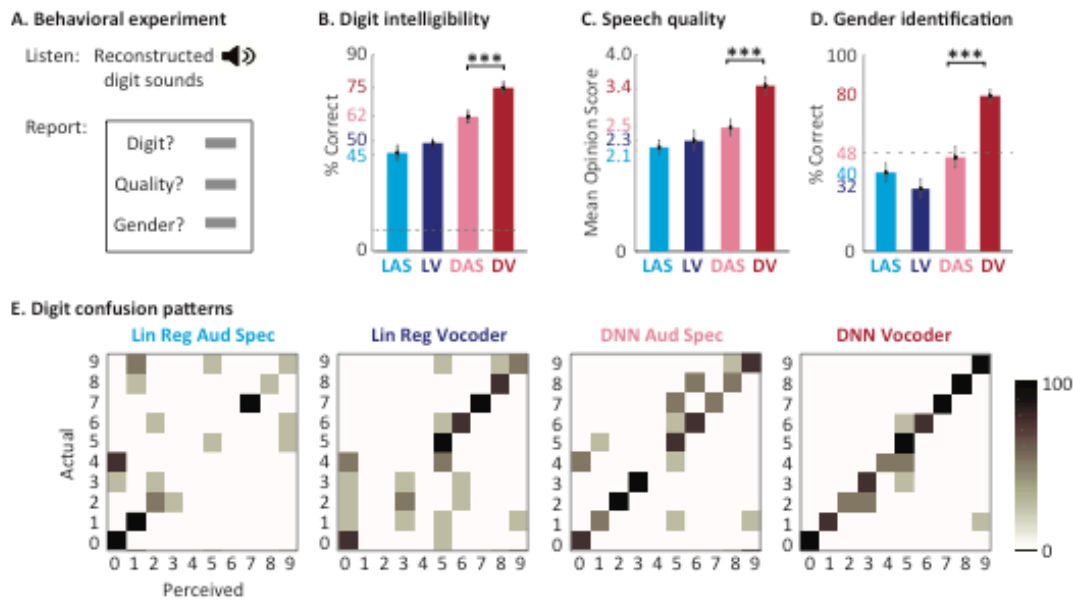
200
201 **Figure 2. Deep neural network architecture** (A) An original auditory spectrogram of a speech sample is
202 shown on top. The reconstructed auditory spectrograms of the four models are shown below. (B)
203 Magnitude power of frequency bands during an unvoiced ($t = 1.4$ sec) and a voiced speech sound ($t = 1.15$
204 sec, shown with dashed lines in A) for original (top) and the four reconstruction models.
205
206

207 **Subjective evaluation of the reconstruction accuracy:** We used the reconstructed digit sounds to assess
208 the subjective intelligibility and quality of the reconstructed audio. Forty unique tokens were
209 reconstructed from each model, consisting of ten digits (zero to nine) that were spoken by two male and
210 two female speakers. The speakers that uttered the digits were different from the speakers that were
211 used in the training, and no digit sound was included in the training of the networks. We asked 11 subjects
212 with normal hearing to listen to the reconstructed digits from all four models (160 tokens total) in a
213 random order. Each digit was heard only once. The subjects then reported the digits (zero to nine, or
214 uncertain), rated the reconstruction quality using the mean opinion score (MOS⁴⁹, on a scale of 1 to 5),
215 and reported the gender of the speaker (Fig. 3A).

216 Figure 3B shows the average reported intelligibility of the digits from the four reconstruction
217 models. The DNN-vocoder combination achieved the best performance (75% accuracy), which is 67%
218 higher than the baseline system (Linear regression with auditory spectrogram). Fig. 3B also shows that the
219 reconstructions using DNN models are significantly better than the linear regression models (68.5% vs.
220 47.5%, paired t-test, $p < 0.001$). Figure 3C shows that the subjects also rated the quality of the
221 reconstruction significantly higher for the DNN-vocoder system than for the other three models (3.4 vs.
222 2.5, 2.3, and 2.1, unpaired t-test, $p < 0.001$), meaning that the DNN-vocoder system sounds closest to
223 natural speech. The subjects also accurately reported the gender of the speaker significantly higher than
224 chance for the DNN-vocoder system (80%, t-test, $p < 0.001$) while the performance for all other methods
225 were at chance (Fig. 3D). The higher intelligibility and quality scores for the DNN-Voc system was
226 consistently observed in all the ten listeners (Supp. Fig. 4). This result indicates the importance of accurate
227 reconstruction of harmonics frequencies for identifying speaker dependent information, which are best
228 captured by the DNN-Voc model.

229 Finally, Figure 3E shows the confusion patterns in recognizing the digits for the four models, confirming
230 again the advantage of the DNN based models, and the DNN vocoder in particular. As shown in Figure 3E,
231 the discriminant acoustic features of the digit sounds are better preserved in the DNN-Voc model,
232 enabling the listeners to correctly differentiate them from the other digits. Linear regression models,
233 however, failed to preserve these cues, as seen by the high confusion among digit sounds. The confusion
234 patterns also show that some errors were associated with the shared phonetic features, for example the
235 confusion between digits one and nine (sharing 'ey' phoneme), or four and five (sharing the initial fricative
236 /f/ phoneme. This result suggests a possible strategy for enabling accurate discrimination in BCI
237 applications by selecting target sounds with a sufficient acoustic distance between them. The audio
238 samples from different models can be found online⁵⁰ and in the supplementary materials.

239



240

241 **Figure 3. Subjective evaluation of the reconstruction accuracy.** (A) The behavioral experiment design

242 used to test the intelligibility and the quality of the reconstructed digits. Eleven subjects listened to digit

243 sounds (zero to nine) spoken by two male and two female speakers. The subjects were asked to report

244 the digit, the quality on the mean-opinion-scale, and the gender of the speaker. (B) The intelligibility score

245 for each model defined as the percentage of correct digits reported by the subject. (C) The quality score

246 on the MOS scale. (D) The speaker gender identification rate for each model. (E) The digit confusion

247 patterns for each of the four models. The DNN vocoder shows the least amount of confusion among the

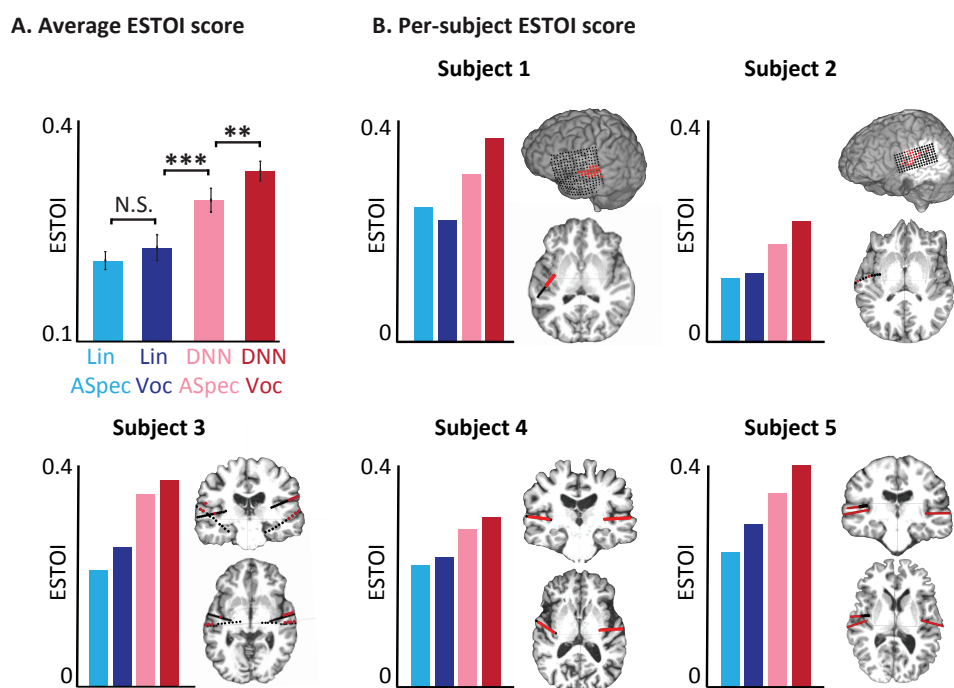
248 digits.

249

250

251

252 **Objective evaluation of reconstructed audio.** We compared the objective reconstruction accuracy of
253 reconstructed audio per subject using the extended short time objective intelligibility (ESTOI)⁵¹ measure.
254 ESTOI is commonly used for the intelligibility assessment of speech synthesis technologies and is
255 calculated by measuring the distortion in spectrotemporal modulation patterns of the noisy speech signal.
256 Therefore, ESTOI score is sensitive to both inaccurate reconstruction of the spectral profile and the
257 inconsistencies in the reconstructed temporal patterns. The ESTOI measures were calculated from
258 continuous speech sentences in the test set. The average ESTOI of the reconstructed speech for all five
259 subjects (Fig. 4A) confirms the results seen from the subjective tests, which is the superiority of DNN based
260 models over the linear model, and that of vocoder reconstruction over the auditory spectrogram
261 ($p < 0.001$, t-test). This pattern was consistent for each of the five subjects in this study, as shown in Fig. 4B
262 alongside the electrode locations for each subject. While the overall reconstruction accuracy varies
263 significantly across subjects, which is likely due to the difference in the coverage of the auditory cortical
264 areas, the relative performance of the four models was the same in all subjects. In addition, averaging the
265 neural responses over multiple repetitions of the same speech utterance improved the reconstruction
266 accuracy (Supp. Fig. 1) because averaging reduces the effect of neural noise.
267



268
269 **Figure 4. Objective intelligibility scores for different models.** (A) The average ESTOI score based on all
270 subjects for the four models. (B) Coverage and the location of the electrodes and ESTOI score for each of
271 the five subjects. In all subjects, the ESTOI score of the DNN vocoder was higher than in the other models.
272
273

274 **Reconstruction accuracy from low and high neural frequencies:** There is increasing evidence that the low
275 and high-frequency bands encode different and complementary information about the stimulus⁵².
276 Considering that the sampling frequency of the reconstruction target is 100 Hz, we used 0–50 Hz as a low-
277 frequency signal, and the envelope of high gamma (70–150 Hz) as high-frequency band information. To
278 determine what frequency bands are best to include to achieve maximum reconstruction accuracy, we
279 tested the reconstruction accuracy in three conditions, when the regression model uses only the high-
280 gamma envelope, a low-frequency signal, or a combination of the two.

281 To simplify the comparison, we used only the DNN-auditory spectrogram reconstruction model.
282 We calculated the ESTOI scores of the reconstructed speech sound using different frequency bands. We
283 found that the combination of the two frequency bands significantly outperforms the reconstruction from
284 only one of the frequency bands (Fig. 5A, $p < 0.001$, t-test). This observation is consistent with the
285 complementary encoding of the stimulus features in the low and high-frequency bands⁵³, which implicates
286 the advantage of using the entire neural signal to achieve the best performance in speech neuroprosthesis
287 applications when it is practically possible.

288
289 **Effect of the number of electrodes and duration of training data:** The variability of the reconstruction
290 accuracy across subjects (Fig. 4B) suggests an important role of neural coverage in improving the
291 reconstruction^{3,6} accuracy. In addition, because some of the noise signal across different electrodes is
292 independent, reconstruction from a combination of electrodes may lead to a higher accuracy by finding a
293 signal subspace less affected by the noise in the data⁵⁴. To examine the effect of the number of electrodes
294 on the reconstruction accuracy, we first combined the electrodes of all five subjects and randomly chose
295 N electrodes ($N = 1, 2, 4, 8, 16, 32, 64, 128$), twenty times for training the individual networks. The average
296 reconstruction accuracy for each N was then used for comparison. The results shown in Fig. 5B indicate
297 that increasing the number of electrodes improves the reconstruction accuracy; however, the rate of
298 improvement decreased significantly.

299 Finally, because the success of neural network models is largely attributed to training on large
300 amounts of data²⁸, we examined the effect of training duration on reconstruction accuracy. We used 128
301 randomly chosen electrodes and trained several neural network models each on a segment of the training
302 data as the duration of the segments was gradually increased from 10 to 30 minutes. This process was
303 performed twenty times for each duration by choosing a random segment of the training data, and the
304 ESTOI score was averaged over the segments. As expected, the results show an increased reconstruction

305 accuracy as the duration of the training was increased (Fig. 5C), which indicates the importance of
306 collecting a larger duration of training data when it is practically feasible.

307

308 **Discussion:**

309 We compared the performance of linear and nonlinear (DNN) regression models in reconstructing the
310 auditory spectrogram and vocoder representation of speech signals. We found that using a deep neural
311 network model to regress vocoder parameters significantly outperformed the linear regression and
312 auditory spectrogram representation of speech, and resulted in 75% intelligibility scores on a closed-set,
313 digit recognition task.

314 Our results are consistent with those of previous reconstruction studies that showed the
315 importance of nonlinear techniques in neural decoding⁵⁵. The previous methods have used support vector
316 machines^{13,56}, linear discriminant analysis^{57,58}, linear regression^{3,14,59}, nonlinear embedding⁶, and Bayes
317 classifiers¹⁵. In recent years, deep learning⁶⁰ has shown tremendous success in many brain-computer
318 interface technologies⁶¹, and our study extended this trend by showing the benefit of deep learning in
319 speech neuroprosthesis research⁵⁵.

320 We showed that the reconstruction accuracy depends on both the number of electrodes and the
321 duration of the data that is available for training. This is consistent with the findings of studies showing
322 the superior advantage of deep learning models over other techniques, particularly when the amount of
323 training data is large²⁸. We showed that the rate of improvement slows down as the number of electrodes
324 increases. This could indicate the limited diversity of the neural responses in our recording which
325 ultimately limits the added information that is gained from additional electrodes. Alternatively, increasing
326 the number of electrodes also increases the complexity and the number of free parameters in the neural
327 network model. Because the duration of our training data was limited, it is possible that more training
328 data would be needed before the benefit of additional features becomes apparent. Our experiments
329 showed that increasing the amount of training data results in better reconstruction accuracy, therefore
330 recording methods that can increase the amount of data available for the training of deep models are
331 highly desirable, for example, when chronic recordings are possible in long-term implantable devices such
332 as the NeuroPace responsive neurostimulation device (RNS)⁶².

333 We showed that the representation of the acoustic signal used as the target of reconstruction has
334 an important role in the intelligibility and the quality of the reconstructed audio. We used a vocoder
335 representation of speech, which extends the previous studies that used a magnitude spectrogram (time-
336 frequency representation)^{3,20}, speech envelope^{21,22}, spectrotemporal modulation frequencies^{6,13,23}, and

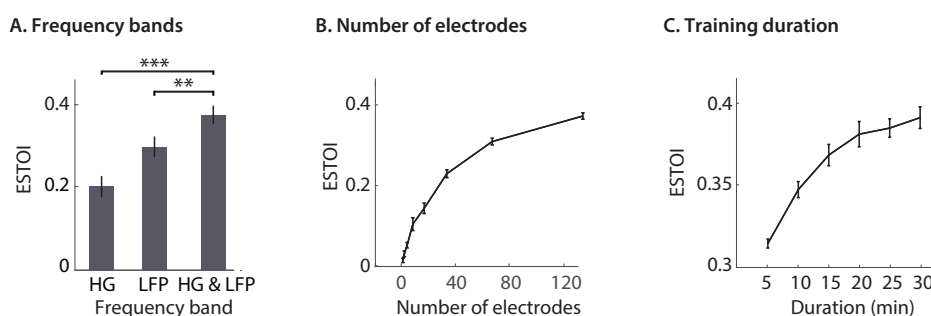
337 discrete units such as phonemes and phonetic categories^{8,17,24,25} and words^{18,19}. Reconstruction of the
338 auditory spectrogram, which we also used for comparison, inherently results in suboptimal audio quality
339 because the phase of the auditory spectrogram must be approximated. The discrete units such as
340 phonemes enable discriminative training by learning a direct map from the neural data to the class labels,
341 which is typically more efficient than generative regression models⁶³. The continuous nature of
342 parameters in acoustic reconstruction however could prove advantageous for BCI applications because
343 they provide a continuous feedback to the user⁶⁴, which is crucial for the subject and the BCI algorithm to
344 coadapt to increase overall effectiveness^{26,27}. Therefore, direct reconstruction of speech synthesis
345 parameters is a natural choice. This choice however poses a challenge, since the vocoder quality is very
346 sensitive to the quality of the decoding. As we have reported, reconstructing vocoder parameters resulted
347 in both the worst (when used with linear regression) and the best (when used with DNN) results.
348 Therefore, powerful modeling techniques such as deep learning are crucial as more inclusive
349 representations of the speech signal are used for reconstruction and decoding applications. We proposed
350 a solution to this problem by compressing the acoustic features into a low-dimensional space and using a
351 decoder that is robust to the fluctuations of the input.

352 We found that the combination of low frequency and the envelope of high gamma results in
353 higher reconstruction accuracy than each frequency band alone. This finding is consistent with those of
354 studies that have shown the importance of an oscillatory phase⁶⁵ in addition to the neural firing rate,
355 which is reflected in the high-gamma frequency band⁶⁶. Combining both high and low frequencies not
356 only enables access to the complementary information in each band^{52,67} but also allows the decoder to
357 use the information that is encoded in the interactions between the two bands, such as cross-frequency
358 coupling⁵³. Overall, we observed that better brain coverage, more training data, and combined neural
359 frequency bands result in the best reconstruction accuracy, which can serve as an upper bound
360 performance where practical limitations prevent the use of all possible factors, for example, where the
361 brain coverage is small, or high-frequency neural signals are not accessible such as in noninvasive
362 neuroimaging methods.

363 The application of neural speech decoding in neuroprosthesis is contingent on the similarity of
364 the underlying neural code in overt and covert (imagined) conditions. Several previous studies have
365 examined the generalization of decoding techniques from overt to covert speech¹²⁻¹⁶ and showed the
366 involvement of the auditory cortical areas, including the superior temporal gyrus (STG) in covert speech
367 condition. Specifically, informative electrodes for speech decoding were found in Wernicke and the STG
368 during imagined articulation^{13,14}, covert word repetition¹⁵, and reading silently¹⁶. In addition to imagined

369 articulation, an MEG study¹² measured the neural activity during actual and imagined hearing conditions
370 and compared with actual and imagined articulation conditions. This study found that the neural activity
371 during overt and covert states were more similar in hearing than in articulation condition. Furthermore,
372 the similarity of the response topographies found in covert and overt hearing suggested a similar neural
373 code in the two states, which is also consistent with the findings of fMRI studies showing a similar neural
374 substrate mediating auditory perception and imagery⁶⁸⁻⁷⁰. It is also worth mentioning that the activation
375 of the auditory cortex is not specific to speech imagery, as a recent study found similar response patterns
376 also during music perception and imagery⁷¹. While these studies have established the feasibility of speech
377 decoding in covert speech perception and production, further research is needed to devise system
378 architectures and training procedures that can optimally fine-tune a model to perform and generalize well
379 in both overt and covert conditions. Furthermore, expanding from the closed-set intelligible speech in this
380 work to continuous, open-set, natural intelligible speech requires additional research, which will
381 undoubtedly benefit from a larger amount of training data, higher-resolution neural recording
382 technologies⁷², and the adaptation of regression models⁷³ and the subject to improve the BCI system^{26,27}.

383 In summary, we present a general framework that can be used for speech neuroprosthesis
384 technologies that can result in accurate and intelligible reconstructed speech from the human auditory
385 cortex. Our approach takes a step toward the next generation of human-computer interaction systems
386 and more natural communication channels for patients suffering from paralysis and locked-in syndromes.



387
388 **Figure 5. Effect of neural frequency range, number of electrodes, and stimulus duration on**
389 **reconstruction accuracy.** (A) The reconstruction ESTOI score based on high gamma, low frequency, and
390 high gamma and low frequency combined. (B) The accuracy of reconstruction when the number of
391 electrodes increases from one to 128. For each condition, 20 random subsets were chosen. (C) The
392 accuracy of reconstruction when the duration of the training data increases. Each condition is the average
393 of 20 random subsets.

394 **Materials and methods:**

395 **Participants and neural recording**

396 Five patients with pharmaco-resistant focal epilepsy were included in this study. All subjects underwent
397 chronic intracranial encephalography (iEEG) monitoring at Northshore University Hospital to identify
398 epileptogenic foci in the brain for later removal. Three subjects were implanted with only stereo-
399 electroencephalographic (sEEG) depth arrays, one with a high-density grid, and one with both grid and
400 depth electrodes (PMT, Chanhassen, MN, USA). The electrodes showing any sign of abnormal epileptiform
401 discharges, as identified in the epileptologists' clinical reports, were excluded from the analysis. All
402 included iEEG time series were manually inspected for signal quality and were free from interictal spikes.
403 All research protocols were approved and monitored by the institutional review board at the Feinstein
404 Institute for Medical Research, and informed written consent to participate in the research studies was
405 obtained from each subject before electrode implantation. All research was performed in accordance with
406 relevant guidelines and regulations.

407 Intracranial EEG (iEEG) signals were acquired continuously at 3 kHz per channel (16-bit precision,
408 range ± 8 mV, DC) using a data acquisition module (Tucker-Davis Technologies, Alachua, FL, USA). Either
409 subdural or skull electrodes were used as references, as dictated by recording quality at the bedside after
410 online visualization of the spectrogram of the signal. Speech signals were recorded simultaneously with
411 the iEEG for subsequent offline analysis. Two ranges of neural frequencies were used in the study. Low-
412 frequency (0–50 Hz) components of the neural data were extracted by filtering the neural signals using an
413 FIR lowpass filter. The high-gamma (70–150 Hz) envelope³⁶ was extracted by first filtering the data into
414 eight frequency bands between 70 and 150 Hz using IIR filters. The envelope of each band was then
415 obtained using a Hilbert transform. We took the average of envelopes in all frequency bands as the total
416 envelope which was then resampled to 100 Hz. The high-gamma responses were normalized based on the
417 responses recorded during a 2-minute silence interval before each recording.

418

419 **Brain maps**

420 The electrode positions were mapped to brain anatomy using registration of the post-implant computed
421 tomography (CT) to the pre-implant MRI via the post-op MRI⁷⁴. After coregistration, the electrodes were
422 identified on the post-implantation CT scan using BioImage Suite⁷⁵. Following coregistration, the subdural
423 grid and strip electrodes were snapped to the closest point on the reconstructed brain surface of the pre-
424 implantation MRI. We used the FreeSurfer automated cortical parcellation⁷⁶ to identify the anatomical
425 regions in which each electrode contact was located within approximately 3 mm resolution (the maximum

426 parcellation error of a given electrode to a parcellated area was < 5 voxels/mm). We used Destrieux's
427 parcellation because it provides higher specificity in the ventral and lateral aspects of the medial lobe⁷⁷.
428 The automated parcellation results for each electrode were closely inspected by a neurosurgeon using the
429 patient's coregistered post-implant MRI.

430

431 **Stimulus**

432 The speech materials included continuous speech stories recorded in-house by four voice actors and
433 actresses (duration: 30 min, 11,025 Hz sampling rate). Eight of the sentences (40 seconds) were used for
434 objective tests and were presented to the patients eight times to improve the signal to noise ratio. The
435 digit sounds were taken from the TI-46 corpus³⁵. Two female (f2 and f8) and two male (m2 and m5)
436 speakers were chosen from the corpus, and one token per digit and speaker was used (total of 40 unique
437 tokens). Each digit was repeated six times to improve the signal to noise ratio of the neural responses.
438 The speakers that uttered the digits were different from the speakers that narrated the stories.

439

440 **Acoustic representation**

441 The auditory spectrogram representation of speech was calculated from a model of the peripheral
442 auditory system⁷⁸. The model consists of three stages: 1) a cochlear filter bank consisting of 128 constant-
443 Q filters equally spaced on a logarithmic axis, 2) a hair cell stage consisting of a low-pass filter and a
444 nonlinear compression function, and 3) a lateral inhibitory network, consisting of a first-order derivative
445 along the spectral axis. Finally, the envelope of each frequency band was calculated to obtain a time-
446 frequency representation simulating the pattern of activity on the auditory nerve⁷⁸. The final auditory
447 spectrogram has a sampling frequency of 100 Hz. The audio signal was reconstructed from the auditory
448 spectrogram using an iterative convex optimization procedure⁴³. For the vocoder-based speech
449 synthesizer, we used the WORLD⁴⁵ (D4C edition) system. In this model, four major speech parameters
450 were estimated, from which the speech waveform was synthesized: 1) spectral envelope, 2) f_0 or
451 fundamental frequency, 3) band aperiodicity, and 4) voiced-unvoiced (VUV) excitation label. The
452 dimension of each parameter was automatically calculated by the vocoder method and was based on the
453 window size and the sampling frequency of the waveform (16 KHz).

454

455 **DNN architecture**

456 We used a common deep neural network architecture that consists of two stages: feature extraction and
457 feature summation³⁸⁻⁴⁰ (Fig. 2A). In this framework, a high-dimensional representation of the input is first

458 calculated (feature extraction), which is then used to regress the output of the model (feature
459 summation). The feature summation and feature extraction networks are optimized jointly together
460 during the training phase. In all models examined, the feature summation step consisted of a two-layer
461 fully connected network (FCN) with L2 regularization, dropout⁴¹, batch normalization⁴², and nonlinearity
462 in each layer.

463 We study five different architectures for the feature extraction part of the network: the fully
464 connected network (FCN, also known as the multilayer perceptron or MLP), the locally connected network
465 (LCN)⁷⁹, convolutional neural network (CNN)⁸⁰, FCN+CNN, and FCN+LCN (for details of each architecture
466 see Supp. Table 1). In the combined networks, we concatenated the output of two parallel paths, which
467 were fed into the summation network. For FCN, the windowed neural responses were flattened and fed
468 to a multilayer FCN. However, in LCN and CNN, all the extracted features were of the same size as the
469 input, meaning that we did not use flattening, strided convolution, or downsampling prior to the input
470 layer or between the two consecutive layers. Instead, the final output of the multilayer LCN or CNN was
471 flattened prior to feeding the output into the feature summation network.

472 The optimal network structure was found separately for the auditory spectrogram and vocoder
473 parameters using an ablation study. For auditory spectrogram reconstruction, we directly regressed the
474 128 frequency bands using a multilayer FCN model for feature extraction (Supp. Table 5). This
475 architecture, however, was not plausible for reconstructing vocoder parameters due to the high-
476 dimensionality and statistical variability of the vocoder parameters. To remedy this, we used a deep
477 autoencoder network (AEC)⁴⁷ to find a compact representation of the 516-dimensional vocoder
478 parameters (consisting of 513 spectral envelopes, pitch, voiced-unvoiced, and band periodicity)⁴⁵. We
479 confirmed that decoding the AEC features performed significantly better than decoding the vocoder
480 parameters directly (Supp. Table 2). The structure for the proposed deep AEC is illustrated in Figure 2D.
481 To carry out decoding, we used a multilayer FCN, in which the number of the nodes changed in a
482 descending (encoder) and then ascending order (decoder) (Fig. 2C)(supp. Table 6). The bottleneck layer
483 of such a network (or the output of the encoder part of the pre-trained AEC) can be used as a low-
484 dimensional reconstruction target by employing the neural network model, from which the vocoder
485 parameters can be estimated using the decoder part of the AEC. We chose the number of nodes in the
486 bottleneck layer to be 256, because it maximized both the objective reconstruction accuracy (Supp. Table
487 3), and the subjective assessment of the reconstructed sound. To increase the robustness to unwanted
488 variations in the encoded features, we used two methods in the bottleneck layer: 1) the hyperbolic
489 tangent function (tanh) was used as a nonlinearity to control the range of the encoded features, and 2)

490 Gaussian noise was added during training prior to feeding into the first layer of the decoder part to make
491 the decoder robust enough to unwanted changes in amplitude resulting from noises in neural responses.
492 We confirmed that using additive Gaussian noise in the bottleneck instead of dropout performed
493 significantly better (paired t-test, $p < 0.001$). It is important that we use the same nonlinearity as the
494 bottleneck (tanh) in the output of the main network, since the estimations should be in the same range
495 and space as those in which they were originally coded. The best network architecture for decoding the
496 vocoder parameters was found to be the FCN+LCN network (Supp. Table 4).

497

498 **DNN training and cross validation**

499 The networks were implemented in Keras with a Tensorflow backend⁸¹. Initialization of the weights was
500 performed using a previously proposed method which was specifically developed for deep multilayer
501 networks with rectified linear units (ReLU) as their nonlinearities⁸². It has been shown that using this
502 method helps such networks converge faster. We used batch normalization⁴², nonlinearity, and a dropout
503 of $p=0.3$ ⁴¹ between each layer. We applied an $L2$ penalty (with a multiplier weight set to 0.001) on the
504 weights of all the layers in all types of networks (including the AEC). However, we found that using additive
505 Gaussian noise in the bottleneck of the AEC instead of dropout and regularization performed significantly
506 better (paired t-test, $p < 0.001$). We used three types of nonlinearities in the networks: 1) LeakyReLU⁸³ for
507 all layers of AEC except the bottleneck and for all layers of the feature extraction part of the main network,
508 2) tanh for the output layer of the main network and the bottleneck of the AEC, and 3) the exponential
509 linear unit (ELU)⁸⁴ for the feature summation network. Each epoch of training had a batch size of 256, and
510 optimization was performed using Adam⁸⁵ with an initial learning rate of 0.0001, which was reduced by a
511 factor of two if the validation loss did not improve in four consecutive epochs. Network training was
512 achieved in 150 epochs and was performed for each subject separately. The loss function was a
513 combination of MSE and Pearson's correlation coefficient for each sample:

514

$$515 \quad \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 - \frac{\sum_i (y_i - \bar{y}_i) (\hat{y}_i - \bar{\hat{y}}_i)}{\sqrt{\sum_i (y_i - \bar{y}_i)^2 \sum_i (\hat{y}_i - \bar{\hat{y}}_i)^2}}$$

516

517 in which y is the actual label (auditory spectrogram or vocoder features) for that sample and \hat{y} is the
518 reconstruction from the output layer of the network. The maximum time-lag used was $\tau_{max} = 300$ ms
519 (Supp. Fig. 2). Because of the higher correlated activity between the neural responses of neighboring
520 electrodes⁸⁶, it was important to ensure that the networks can model the local structure in the data.

521 Because both CNN and LCN use small receptive fields that take local patterns into account, we retained
522 the spatial organization of the electrode sites in the input to the network, meaning that the electrodes
523 that were close to each other in the brain were arranged to be close together in the input data matrix.

524

525 **Cross validation**

526 We trained both the LR model and the DNN models using cross validation. We used the speech stories for
527 training all models, and used repeated sentences (separate set from the stories) and digit sounds for
528 testing. No digit sound was included in the training, and the speakers that uttered the digits were different
529 from those that read the stories. The autoencoder network (AEC) was trained on a separate speech corpus
530 (Wall Street Journal, WSJ, 80 hours of read speech)⁴⁸.

531

532 **Subjective and objective evaluations**

533 We assessed the intelligibility of the reconstructed speech using both subjective and objective tests. For
534 subjective assessment, 11 participants with self-reported normal hearing listened to the reconstructed
535 digits using headphones in a quiet environment. Each participant listened to 160 tokens including 10
536 digits, four speakers, and four models. The participants were asked to report the digit or to select unsure
537 if the digit was not intelligible. In addition, the participants reported the quality of the reconstructed
538 speech using a mean opinion score (MOS): 1 (bad), 2 (poor), 3 (fair), 4 (good), and 5 (excellent). The
539 participants also reported the gender of the speaker. For objective evaluation, we used the ESTOI
540 measure⁵¹ which is a monaural intelligibility prediction algorithm commonly used in speech enhancement
541 and synthesis research. The range of ESTOI measure is between zero (worst) and one (best).

542

543

544 **Data availability**

545 The data that support the findings of this study are available upon request from the corresponding author
546 [NM].

547

548 **Code availability**

549 The codes for performing phoneme analysis, calculating high-gamma envelope, and reconstructing the
550 auditory spectrogram are available at <http://naplab.ee.columbia.edu/naplib.html>⁸⁷.

551

552 **Acknowledgments**

553 We thank James O’Sullivan for providing helpful comments on the manuscript. This work was funded by
554 a grant from the National Institutes of Health, NIDCD, DC014279 and the Pew Charitable Trusts, Pew
555 Biomedical Scholars Program.

556

557 **Author Contributions**

558 H.A., B.K., N.M., designed the experiment, evaluated the results and wrote the manuscript. B.K., J.H., N.M.,
559 A.M. collected the data. All authors commented on the manuscript.

560

561 **Competing interests**

562 The authors declare no competing interests.

563

564

565

566 **References**

567

- 568 1. Bialek, W., Rieke, F., de Ruyter van Steveninck, R. R. & Warland, D. Reading a neural code. *Science*
569 (80-.). **252**, 1854–1857 (1991).
- 570 2. Rieke, F., Bodnar, D. A. & Bialek, W. Naturalistic stimuli increase the rate and efficiency of
571 information transmission by primary auditory afferents. *Proc Biol Sci* **262**, 259–265 (1995).
- 572 3. Mesgarani, N., David, S. V. S. V., Fritz, J. B. J. B. & Shamma, S. A. S. A. Influence of context and
573 behavior on stimulus reconstruction from neural activity in primary auditory cortex. *J*
574 *Neurophysiol* **102**, 3329–3339 (2009).
- 575 4. Stanley, G. B., Li, F. F. & Dan, Y. Reconstruction of natural scenes from ensemble responses in the
576 lateral geniculate nucleus. *J Neurosci* **19**, 8036–8042 (1999).
- 577 5. Ramirez, A. D. *et al.* Incorporating naturalistic correlation structure improves spectrogram
578 reconstruction from neuronal activity in the songbird auditory midbrain. *J. Neurosci.* **31**, 3828–
579 3842 (2011).
- 580 6. Pasley, B. N. B. N. *et al.* Reconstructing speech from human auditory cortex. *PLoS Biol.* **10**, (2012).
- 581 7. Chakrabarti, S., Sandberg, H. M., Brumberg, J. S. & Krusienski, D. J. Progress in speech decoding
582 from the electrocorticogram. *Biomed. Eng. Lett.* **5**, 10–21 (2015).
- 583 8. Herff, C. & Schultz, T. Automatic speech recognition from neural signals: a focused review. *Front.*
584 *Neurosci.* **10**, 429 (2016).
- 585 9. Iljina, O. *et al.* Neurolinguistic and machine-learning perspectives on direct speech BCIs for
586 restoration of naturalistic communication. *Brain-Computer Interfaces* **4**, 186–199 (2017).
- 587 10. Laureys, S. *et al.* The locked-in syndrome: what is it like to be conscious but paralyzed and
588 voiceless? *Prog. Brain Res.* **150**, 495–611 (2005).
- 589 11. Sellers, E. W., Ryan, D. B. & Hauser, C. K. Noninvasive brain-computer interface enables
590 communication after brainstem stroke. *Sci. Transl. Med.* **6**, 257re7-257re7 (2014).
- 591 12. Tian, X. & Poeppel, D. Mental imagery of speech and movement implicates the dynamics of
592 internal forward models. *Front. Psychol.* **1**, 166 (2010).
- 593 13. Martin, S. *et al.* Word pair classification during imagined speech using direct brain recordings. *Sci.*
594 *Rep.* **6**, 25803 (2016).
- 595 14. Leuthardt, E. C. *et al.* Using the electrocorticographic speech network to control a brain–
596 computer interface in humans. *J. Neural Eng.* **8**, 36004 (2011).
- 597 15. Pei, X., Barbour, D. L., Leuthardt, E. C. & Schalk, G. Decoding vowels and consonants in spoken

- 598 and imagined words using electrocorticographic signals in humans. *J. Neural Eng.* **8**, 46028
599 (2011).
- 600 16. Martin, S. *et al.* Decoding spectrotemporal features of overt and covert speech from the human
601 cortex. *Front. Neuroeng.* **7**, 14 (2014).
- 602 17. Di Liberto, G. M., O'Sullivan, J. A. & Lalor, E. C. Low-Frequency Cortical Entrainment to Speech
603 Reflects Phoneme-Level Processing. *Curr. Biol.* **25**, 2457–2465 (2015).
- 604 18. Kellis, S. *et al.* Decoding spoken words using local field potentials recorded from the cortical
605 surface. *J. Neural Eng.* **7**, 56007 (2010).
- 606 19. Herff, C. *et al.* Brain-to-text: decoding spoken phrases from phone representations in the brain.
607 *Front. Neurosci.* **9**, 217 (2015).
- 608 20. Mesgarani, N. & Chang, E. F. E. F. Selective cortical representation of attended speaker in multi-
609 talker speech perception. *Nature* **485**, 233–236 (2012).
- 610 21. O'Sullivan, J. A. *et al.* Attentional Selection in a Cocktail Party Environment Can Be Decoded from
611 Single-Trial EEG. *Cereb. Cortex* bht355 (2014).
- 612 22. Ding, N. & Simon, J. Z. Emergence of neural encoding of auditory objects while listening to
613 competing speakers. *Proc. Natl. Acad. Sci.* **109**, 11854–11859 (2012).
- 614 23. Santoro, R. *et al.* Reconstructing the spectrotemporal modulations of real-life sounds from fMRI
615 response patterns. *Proc. Natl. Acad. Sci.* **114**, 4799–4804 (2017).
- 616 24. Moses, D. A. D. A., Mesgarani, N., Leonard, M. K. M. K. & Chang, E. F. E. F. Neural speech
617 recognition: continuous phoneme decoding using spatiotemporal representations of human
618 cortical activity. *J. Neural Eng.* **13**, 56004 (2016).
- 619 25. Khalighinejad, B., da Silva, G. C. & Mesgarani, N. Dynamic Encoding of Acoustic Features in Neural
620 Responses to Continuous Speech. *J. Neurosci.* **37**, 2176–2185 (2017).
- 621 26. Vidaurre, C., Sannelli, C., Müller, K.-R. & Blankertz, B. Machine-learning-based coadaptive
622 calibration for brain-computer interfaces. *Neural Comput.* **23**, 791–816 (2011).
- 623 27. McFarland, D. J., Sarnacki, W. A. & Wolpaw, J. R. Should the parameters of a BCI translation
624 algorithm be continually adapted? *J. Neurosci. Methods* **199**, 103–107 (2011).
- 625 28. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436 (2015).
- 626 29. Hinton, G. *et al.* Deep neural networks for acoustic modeling in speech recognition: The shared
627 views of four research groups. *Signal Process. Mag. IEEE* **29**, 82–97 (2012).
- 628 30. Luo, Y. Y., Chen, Z. & Mesgarani, N. Speaker-Independent Speech Separation With Deep Attractor
629 Network. *IEEE/ACM Trans. Audio, Speech, Lang. Process.* **26**, 787–796 (2018).

- 630 31. Chen, Z., Luo, Y. Y. & Mesgarani, N. Deep attractor network for single-microphone speaker
631 separation. in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International*
632 *Conference on* 246–250 (IEEE, 2017). doi:10.1109/ICASSP.2017.7952155
- 633 32. O’Sullivan, J. *et al.* Neural decoding of attentional selection in multi-speaker environments
634 without access to clean sources. *J. Neural Eng.* **14**, 56001 (2017).
- 635 33. David, S. V. S. V., Mesgarani, N., Fritz, J. B. J. B. & Shamma, S. A. S. A. Rapid synaptic depression
636 explains nonlinear modulation of spectro-temporal tuning in primary auditory cortex by natural
637 stimuli. *J Neurosci* **29**, 3374–3386 (2009).
- 638 34. Mesgarani, N., David, S. V. S. V., Fritz, J. B. J. B. & Shamma, S. A. S. A. Mechanisms of noise robust
639 representation of speech in primary auditory cortex. *Proc. Natl. Acad. Sci.* **111**, 6792–6797
640 (2014).
- 641 35. Liberman, Mark, *et al.* TI 46-Word LDC93S9. *Linguistic Data Consortium, Philadelphia* (1993).
- 642 36. Crone, N. E., Boatman, D., Gordon, B. & Hao, L. Induced electrocorticographic gamma activity
643 during auditory perception. *Clin. Neurophysiol.* **112**, 565–582 (2001).
- 644 37. Edwards, E. *et al.* Comparison of time–frequency responses and the event-related potential to
645 auditory speech stimuli in human cortex. *J. Neurophysiol.* **102**, 377–386 (2009).
- 646 38. LeCun, Y. *et al.* Handwritten digit recognition with a back-propagation network. in *Advances in*
647 *neural information processing systems* 396–404 (1990).
- 648 39. Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural
649 networks. in *Advances in neural information processing systems* 1097–1105 (2012).
- 650 40. Pinto, N., Doukhan, D., DiCarlo, J. J. & Cox, D. D. A high-throughput screening approach to
651 discovering good forms of biologically inspired visual representation. *PLoS Comput. Biol.* **5**,
652 e1000579 (2009).
- 653 41. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A simple way
654 to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
- 655 42. Ioffe, S. & Szegedy, C. Batch normalization: Accelerating deep network training by reducing
656 internal covariate shift. *arXiv Prepr. arXiv1502.03167* (2015).
- 657 43. Chi, T., Ru, P. & Shamma, S. A. Multiresolution spectrotemporal analysis of complex sounds. *J*
658 *Acoust Soc Am* **118**, 887–906 (2005).
- 659 44. Mesgarani, N., Slaney, M. & Shamma, S. A. S. A. Discrimination of speech from nonspeech based
660 on multiscale spectro-temporal modulations. *IEEE Trans. Audio. Speech. Lang. Processing* **14**,
661 920–930 (2006).

- 662 45. Morise, M., Yokomori, F. & Ozawa, K. WORLD: a vocoder-based high-quality speech synthesis
663 system for real-time applications. *IEICE Trans. Inf. Syst.* **99**, 1877–1884 (2016).
- 664 46. Kawahara, H., Masuda-Katsuse, I. & De Cheveigne, A. Restructuring speech representations using
665 a pitch-adaptive time–frequency smoothing and an instantaneous-frequency-based F0
666 extraction: Possible role of a repetitive structure in sounds¹. *Speech Commun.* **27**, 187–207
667 (1999).
- 668 47. Hinton, G. E. & Salakhutdinov, R. R. Reducing the dimensionality of data with neural networks.
669 *Science (80-.)*. **313**, 504–507 (2006).
- 670 48. Paul, D. B. & Baker, J. M. The design for the Wall Street Journal-based CSR corpus. in *Proceedings*
671 *of the workshop on Speech and Natural Language* 357–362 (Association for Computational
672 Linguistics, 1992).
- 673 49. Salza, P. L., Foti, E., Nebbia, L. & Oreglia, M. MOS and pair comparison combined methods for
674 quality evaluation of text-to-speech systems. *Acta Acust. united with Acust.* **82**, 650–656 (1996).
- 675 50. Reconstruction audio samples: naplab.columbia.edu/Reconstruction.
- 676 51. Jensen, J. & Taal, C. H. An Algorithm for Predicting the Intelligibility of Speech Masked by
677 Modulated Noise Maskers. *IEEE/ACM Trans. Audio, Speech Lang. Process.* **24**, 2009–2022 (2016).
- 678 52. Buzsáki, G., Anastassiou, C. A. & Koch, C. The origin of extracellular fields and currents—EEG,
679 ECoG, LFP and spikes. *Nat. Rev. Neurosci.* **13**, 407–420 (2012).
- 680 53. Canolty, R. T. & Knight, R. T. The functional role of cross-frequency coupling. *Trends Cogn. Sci.* **14**,
681 506–515 (2010).
- 682 54. Paninski, L., Pillow, J. & Lewi, J. Statistical models for neural encoding, decoding, and optimal
683 stimulus design. *Prog. Brain Res.* **165**, 493–507 (2007).
- 684 55. Yang, M. *et al.* Speech reconstruction from human auditory cortex with deep neural networks. in
685 *Sixteenth Annual Conference of the International Speech Communication Association* (2015).
- 686 56. Blakely, T., Miller, K. J., Rao, R. P. N., Holmes, M. D. & Ojemann, J. G. Localization and
687 classification of phonemes using high spatial resolution electrocorticography (ECoG) grids. *Conf.*
688 *Proc. IEEE Eng. Med. Biol. Soc.* **2008**, 4964–7 (2008).
- 689 57. Mugler, E. M. *et al.* Direct classification of all American English phonemes using signals from
690 functional speech motor cortex. *J. Neural Eng.* **11**, 35015 (2014).
- 691 58. Lotte, F. *et al.* Electrocorticographic representations of segmental features in continuous speech.
692 *Front. Hum. Neurosci.* **9**, 97 (2015).
- 693 59. Herff, C. *et al.* Towards direct speech synthesis from ECoG: A pilot study. in *Engineering in*

- 694 *Medicine and Biology Society (EMBC), 2016 IEEE 38th Annual International Conference of the*
695 *1540–1543 (IEEE, 2016).*
- 696 60. Hinton, G. E. *et al.* A fast learning algorithm for deep belief nets. *Neural Comput.* **18**, 1527–1554
697 (2006).
- 698 61. Hajinoroozi, M., Mao, Z., Jung, T.-P., Lin, C.-T. & Huang, Y. EEG-based prediction of driver’s
699 cognitive performance by deep convolutional neural network. *Signal Process. Image Commun.*
700 **47**, 549–555 (2016).
- 701 62. Morrell, M. Brain stimulation for epilepsy: can scheduled or responsive neurostimulation stop
702 seizures? *Curr. Opin. Neurol.* **19**, 164–168 (2006).
- 703 63. Efron, B. The efficiency of logistic regression compared to normal discriminant analysis. *J. Am.*
704 *Stat. Assoc.* **70**, 892–898 (1975).
- 705 64. Koyama, S. *et al.* Comparison of brain–computer interface decoding algorithms in open-loop and
706 closed-loop control. *J. Comput. Neurosci.* **29**, 73–87 (2010).
- 707 65. Luo, H. & Poeppel, D. Phase patterns of neuronal responses reliably discriminate speech in
708 human auditory cortex. *Neuron* **54**, 1001–1010 (2007).
- 709 66. Ray, S. & Maunsell, J. H. R. Different Origins of Gamma Rhythm and High-Gamma Activity in
710 Macaque Visual Cortex. *PLoS Biol.* **9**, (2011).
- 711 67. Miller, K. J., Sorensen, L. B., Ojemann, J. G. & Den Nijs, M. Power-law scaling in the brain surface
712 electric potential. *PLoS Comput. Biol.* **5**, e1000609 (2009).
- 713 68. Zatorre, R. J., Halpern, A. R., Perry, D. W., Meyer, E. & Evans, A. C. Hearing in the mind’s ear: a
714 PET investigation of musical imagery and perception. *J. Cogn. Neurosci.* **8**, 29–46 (1996).
- 715 69. Aleman, A. *et al.* The functional neuroanatomy of metrical stress evaluation of perceived and
716 imagined spoken words. *Cereb. Cortex* **15**, 221–228 (2005).
- 717 70. Bunzeck, N., Wuestenberg, T., Lutz, K., Heinze, H.-J. & Jancke, L. Scanning silence: mental imagery
718 of complex sounds. *Neuroimage* **26**, 1119–1127 (2005).
- 719 71. Martin, S. *et al.* Neural encoding of auditory features during music perception and imagery.
720 *Cereb. Cortex* 1–12 (2017).
- 721 72. Khodagholy, D. *et al.* NeuroGrid: recording action potentials from the surface of the brain. *Nat.*
722 *Neurosci.* **18**, 310 (2015).
- 723 73. Shenoy, P., Krauledat, M., Blankertz, B., Rao, R. P. N. & Müller, K.-R. Towards adaptive
724 classification for BCI. *J. Neural Eng.* **3**, R13 (2006).
- 725 74. Groppe, D. M. *et al.* iELVis: An open source MATLAB toolbox for localizing and visualizing human

- 726 intracranial electrode data. **281**,
- 727 75. Papademetris, X. *et al.* BioImage Suite: An integrated medical image analysis suite: An update.
728 *Insight J.* **2006**, 209 (2006).
- 729 76. Fischl, B. *et al.* Automatically parcellating the human cerebral cortex. *Cereb. cortex* **14**, 11–22
730 (2004).
- 731 77. Destrieux, C., Fischl, B., Dale, A. & Halgren, E. Automatic parcellation of human cortical gyri and
732 sulci using standard anatomical nomenclature. *Neuroimage* **53**, 1–15 (2010).
- 733 78. Yang X. Shamma S. A., W. K. Auditory representations of acoustic signals. *IEEE Trans. Inf. Theory*
734 **38**, 824–839 (1992).
- 735 79. Coates, A. & Ng, A. Y. Selecting receptive fields in deep networks. in *Advances in Neural*
736 *Information Processing Systems* 2528–2536 (2011).
- 737 80. LeCun, Y. & Bengio, Y. Convolutional networks for images, speech, and time series. *Handb. brain*
738 *theory neural networks* **3361**, (1995).
- 739 81. Abadi, M. *et al.* TensorFlow: A System for Large-Scale Machine Learning. in *OSDI* **16**, 265–283
740 (2016).
- 741 82. He, K., Zhang, X., Ren, S. & Sun, J. Delving deep into rectifiers: Surpassing human-level
742 performance on imagenet classification. in *Proceedings of the IEEE international conference on*
743 *computer vision* 1026–1034 (2015).
- 744 83. Maas, A. L., Hannun, A. Y. & Ng, A. Y. Rectifier nonlinearities improve neural network acoustic
745 models. in *Proc. icml* **30**, 3 (2013).
- 746 84. Clevert, D.-A., Unterthiner, T. & Hochreiter, S. Fast and accurate deep network learning by
747 exponential linear units (elus). *arXiv Prepr. arXiv1511.07289* (2015).
- 748 85. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv Prepr. arXiv1412.6980*
749 (2014).
- 750 86. Muller, L., Hamilton, L. S., Edwards, E., Bouchard, K. E. & Chang, E. F. Spatial resolution
751 dependence on spectral frequency in human speech cortex electrocorticography. *J. Neural Eng.*
752 **13**, 56013 (2016).
- 753 87. Khalighinejad, B., Nagamine, T., Mehta, A. & Mesgarani, N. NAPLib: An open source toolbox for
754 real-time and offline Neural Acoustic Processing. in *Acoustics, Speech and Signal Processing*
755 *(ICASSP), 2017 IEEE International Conference on* 846–850 (IEEE, 2017).
756 doi:10.1109/ICASSP.2017.7952275
757