

1 **Parallel hierarchical encoding of linguistic representations in the human**
2 **auditory cortex and recurrent automatic speech recognition systems**

3 Menoua Keshishian ^{1,2}, Gavin Mischler ^{1,2}, Samuel Thomas ³, Brian Kingsbury ³,
4 Stephan Bickel ^{4,5}, Ashesh D. Mehta ^{4,5}, Nima Mesgarani ^{1,2,*}

5 **Abstract**

6 The human brain’s ability to transform acoustic speech signals into rich linguistic representations has
7 inspired advancements in automatic speech recognition (ASR) systems. While ASR systems now achieve
8 human-level performance under controlled conditions, prior research on their parallels with the brain has
9 been limited by the use of biologically implausible models, narrow feature sets, and comparisons that
10 primarily emphasize predictability of brain activity without fully exploring shared underlying
11 representations. Additionally, studies comparing the brain to text-based language models overlook the
12 acoustic stages of speech processing, an essential part in transforming sound to meaning. Leveraging high-
13 resolution intracranial recordings and a recurrent ASR model, this study bridges these gaps by uncovering
14 a striking correspondence in the hierarchical encoding of linguistic features, from low-level acoustic
15 signals to high-level semantic processing. Specifically, we demonstrate that neural activity in distinct
16 regions of the auditory cortex aligns with representations in corresponding layers of the ASR model and,
17 crucially, that both systems encode similar features at each stage of processing—from acoustic to
18 phonetic, lexical, and semantic information. These findings suggest that both systems, despite their
19 distinct architectures, converge on similar strategies for language processing, providing insight in the
20 optimal computational principles underlying linguistic representation and the shared constraints shaping
21 human and artificial speech processing.

22

¹ Department of Electrical Engineering, Columbia University, New York, NY, USA

² Zuckerman Institute, Columbia University, New York, NY, USA

³ IBM Research, Yorktown Heights, NY, USA

⁴ The Feinstein Institutes for Medical Research, Northwell Health, Manhasset, NY, USA

⁵ Department of Neurosurgery, Zucker School of Medicine at Hofstra/Northwell, Hempstead, NY, USA

* Correspondence: nima@ee.columbia.edu

23 Introduction

24 Understanding how the human brain processes and encodes linguistic information is a
25 fundamental challenge in neuroscience and artificial intelligence. The human auditory cortex is
26 capable of extracting meaning and structure from spoken language with remarkable efficiency.
27 In parallel, automatic speech recognition (ASR) systems have achieved near-human accuracy in
28 recognizing and transcribing speech (1). However, it remains unclear how closely the internal
29 computation and representation of these ASR systems mirror those of the human brain, leaving
30 a significant gap in our understanding of speech processing in humans and machines. Are both
31 systems converging toward similar strategies, or are they reaching different solutions
32 independently as they optimize for performance?

33 Neuroimaging studies of speech processing in the brain have revealed an emergent encoding of
34 linguistic hierarchies, progressing from primary to nonprimary areas of the auditory cortex (2–4).
35 These studies show that the brain distributes, yet jointly encodes, various linguistic features,
36 including phonemes, phonotactics, lexical-phonological, and lexical-semantic information.
37 However, parallels between these hierarchical patterns in the brain and those found in ASR
38 algorithms have not been directly established. Several studies have examined how end-to-end
39 ASR systems represent linguistic information (5) and explored similarities between these systems
40 and the brain (6). Additional research has also investigated the representational similarities
41 between large language models (LLMs) and the brain (7–11). While these studies have provided
42 valuable insights, they have some limitations. Many lacked the temporal precision necessary for
43 speech processing, as they relied on functional MRI (fMRI) (7, 10) which is slow in capturing the
44 neural dynamics. Moreover, by comparing the brain to text-based LLMs, these studies inherently
45 overlooked the acoustic stages of speech processing even though the subjects listened to the
46 stimuli (7–11). Other studies used models that were biologically implausible, such as
47 transformers (6–11) or non-causal architectures (6). In addition, some studies did not explicitly
48 analyze linguistic representation in the models (7, 9–11), or used only a narrow set of features
49 (6, 8). Therefore, even though it has been shown that speech processing models and the brain
50 use increasingly similar representations as revealed by the predictability of neural responses

51 from these models, the precise nature of this convergent similarity remains unclear, replacing
52 one black box with another one without explaining the full picture of speech understanding.

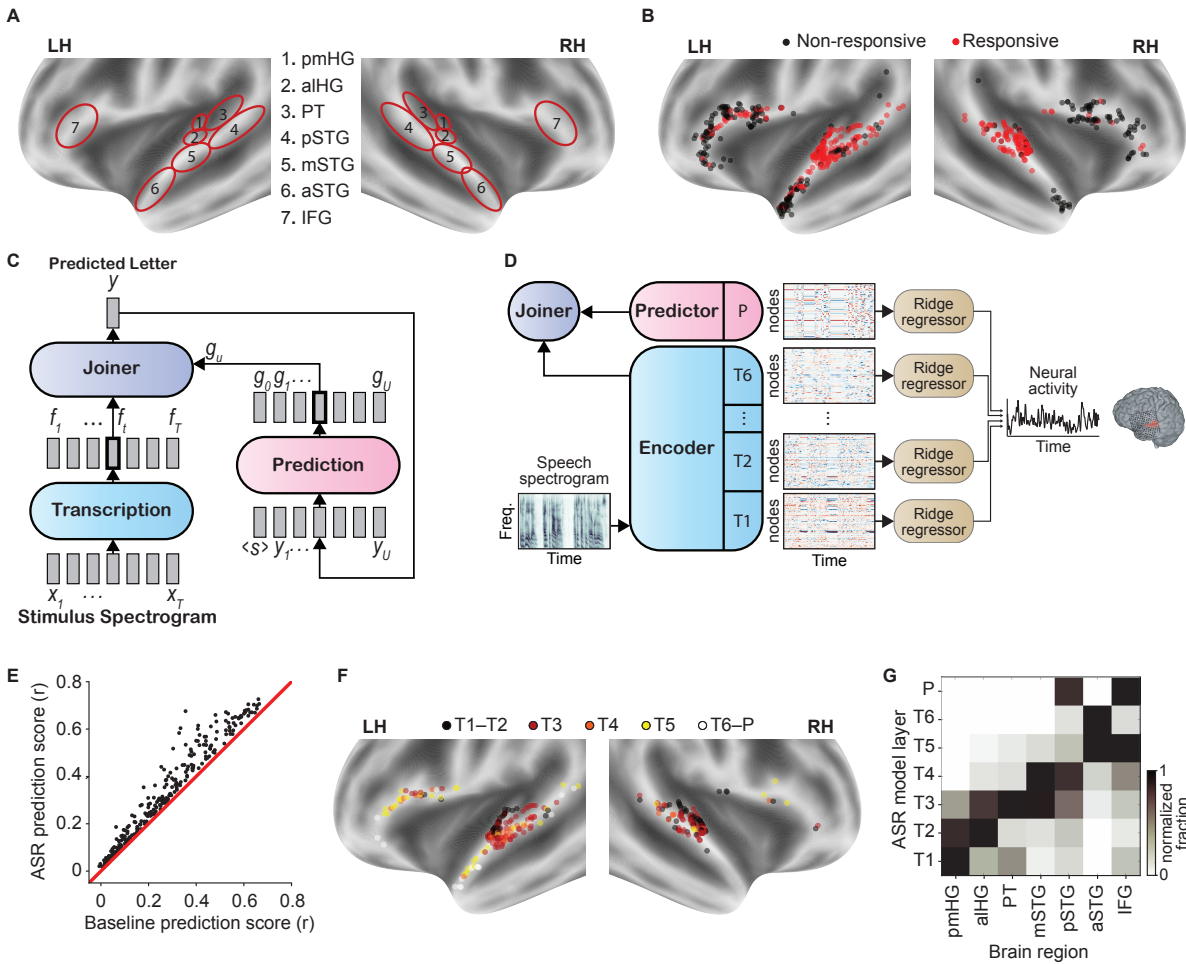
53 This study addresses several gaps in the current understanding of speech processing in both
54 biological and artificial systems. First, we expand on previous studies by investigating a broad
55 range of linguistic features, providing a comprehensive analysis that includes phonetic, lexical,
56 and semantic representations. Second, we use a biologically plausible recurrent neural network
57 transducer (RNN-T) (12) model that processes speech in a causal and incremental manner,
58 aligning more closely with how the human brain processes speech in real time. Finally, by
59 incorporating high-resolution intracranial electroencephalography (iEEG) data from participants
60 listening to continuous speech, we establish a direct comparison between neural activations and
61 the internal states of the ASR model.

62 Through a detailed node-level analysis, we compare specific neural sites in the brain with
63 individual nodes in the ASR model to directly assess representational alignment. Additionally, a
64 layer-level analysis provides a more comprehensive view of the ASR model's internal hierarchical
65 representations, allowing us to examine how the structure of these representations parallels the
66 cortical encoding of speech in humans. This dual approach enables us to uncover both fine-
67 grained and hierarchical similarities in how speech is processed by biological and artificial
68 systems, shedding light on shared mechanisms of linguistic encoding while also revealing key
69 divergences.

70 **Results**

71 We recorded intracranial electroencephalography (iEEG) data from fifteen human participants
72 implanted with subdural (electrocorticography; ECoG) and depth (stereotactic EEG) electrodes.
73 The participants listened to 30 minutes of continuous speech spoken by four speakers. To ensure
74 that the subjects were engaged in the task, we paused the audio at random intervals and asked
75 the subjects to report the last sentence before the pause. All subjects were attentive and could
76 correctly repeat the speech utterances. We extracted the envelope of the high-gamma frequency
77 band (70-150 Hz), which has been shown to correlate with neural firing in the proximity of the
78 recording electrode (13, 14), as the neural response measure of the recorded signals.

79 We restricted our analyses to electrode sites in the auditory cortex (AC; $N = 335$) and the
80 inferior frontal gyrus (IFG; $N = 191$). Figure 1A shows the general location of the IFG and the
81 subregions of the auditory cortex on the FreeSurfer average brain (15). We further limited our
82 analysis to sites that were determined to be speech-responsive ($N = 291/526$, determined by
83 a t-test between responses during speech and in silence). We labeled the neural sites in both
84 hemispheres based on their anatomical region to enable population tests: posteromedial
85 Heschl's gyrus (pmHG; $N_L = 12/15$, $N_R = 16/19$), anterolateral Heschl's gyrus (alHG; $N_L =$
86 $34/36$, $N_R = 32/34$), planum temporale (PT; $N_L = 9/12$, $N_R = 18/20$), middle superior
87 temporal gyrus (mSTG; $N_L = 44/53$, $N_R = 17/18$), posterior superior temporal gyrus (pSTG;
88 $N_L = 27/32$, $N_R = 21/23$), anterior superior temporal gyrus (aSTG; $N_L = 19/60$, $N_R = 0/13$),
89 and inferior frontal gyrus (IFG; $N_L = 32/119$, $N_R = 10/72$). The electrode locations and their
90 responsiveness are plotted in Figure 1B on the average FreeSurfer brain, where the color
91 indicates whether an electrode is speech-responsive.



92

93 **Figure 1. Stages of speech processing in the brain based on ASR modelling.**

94 **(A)** General location of the inferior frontal gyrus (IFG) and subregions of the auditory cortex (AC).

95 **(B)** Electrode locations within the region of interest. Colors represent whether the electrode

96 showed higher activation in the presence of speech stimuli compared to moments of pre-trial

97 silence. **(C)** RNN-Transducer architecture. x_t are spectrogram frames, f_t are the encoder's

98 outputs, $h_{t,u}$ are model outputs used to generate output labels, y_u are output labels predicted

99 by the model (feedback), and g_u are the prediction network's outputs. **(D)** Predicting brain

100 responses from the ASR model. To determine the mapping between network layers and

101 electrode sites, we fit a Ridge regression model to predict the neural response at the electrode

102 (r_S) while the participant was listening to some stimuli S , from the activation matrix $\Gamma_{i,S}$

103 corresponding to the response of the i -th layer of the ASR model to the same stimuli S . **(E)**

104 Improvement in prediction score of electrode sites when using the best ASR layer compared to a

105 spectrogram baseline predictor. The best ASR layer for an electrode was determined by 10-fold
106 cross-validation. **(F)** We color each electrode based on the layer of the ASR model that predicts
107 it best. **(G)** Within each subregion of interest, we compute the fraction of electrodes that are best
108 predicted by each layer of the ASR. We then normalize these fractions by dividing each row by
109 its maximum value. This metric shows in which subregion of interest a given layer of the model
110 is represented relatively stronger.

111

112

113 **Stages of speech processing in the brain based on ASR modeling**

114 We use an RNN-Transducer (RNN-T) (12), a recurrent neural network trained to predict letters
115 (graphemes) from the speech spectrogram, to model the speech recognition process in the brain.
116 RNN-T models take spectro-temporal speech signals as input through a causal mechanism, and
117 have been widely used for speech recognition (16–20), making them an ideal candidate model to
118 compare with the full hierarchy of speech processing in the brain. The model consists of two
119 branches—transcription and prediction (Figure 1C). The transcription branch, consisting of 6 uni-
120 directional LSTM layers (T1–T6), acts as a bottom-up encoder by producing acoustic embeddings
121 from the input speech spectrogram. The prediction branch, consisting of 1 uni-directional LSTM
122 layer (P), acts as an internal language model in that it is conditioned on previous non-blank
123 symbols produced by the model. It stores an internal representation of the prior predictions of
124 the model that together with the encoder’s output are combined by a shallow joint network and
125 used to make the next prediction. We computed layer activations for all layers of these two
126 branches in response to the same 30-minute stimulus set that the human participants listened
127 to.

128 To examine how the speech processing pathways of the brain and the model map to each other,
129 we find the best matching model layer for each electrode site (Figure 1D). First, for each
130 electrode-layer pair, we fit a single-lag regression model that predicts the electrode activity from
131 the layer activations in response to the same stimuli (reduced to top 256 principal components):

$$132 \quad y(t) = X(t - \delta)\beta + c + \varepsilon$$

133 where y is the $T \times 1$ response vector for an electrode site, X is the $T \times H$ matrix of layer
134 activations, δ is a constant time lag, β is the $H \times 1$ vector of regression coefficients, c is a
135 constant bias term and ε is the regression error. We also define a baseline 256-dimensional
136 predictor which is obtained from a 128-dimensional log-Mel-spectrogram stacked on top of a 1-
137 sample delayed version of itself, such that $X(t, 1:128) = S(t, :)$ and $X(t, 129:256) =$
138 $S(t - 1, :)$, where X is the $T \times 256$ baseline predictor and S is the $T \times 128$ spectrogram. For
139 each electrode, we also fit a single-lag linear regression model that predicts the neural activity
140 from the baseline predictor, allowing for a fixed time lag δ . Figure 1E shows the improvement of
141 cross-validated prediction scores for each electrode when predicting from the ASR model
142 compared to the baseline spectrogram predictor.

143 We associate each electrode with the layer of the model that predicts it with the highest
144 accuracy. The colors in Figure 1F indicate the corresponding layer of the model for each neural
145 site in the AC and IFG; and Figure 1G shows the normalized fraction of electrodes in each brain
146 region best predicted by each layer of the model (values indicate fraction of neural sites in brain
147 region, divided by sum of fractions within layer/row). The model layers map to the cortex such
148 that as we move deeper in the model, in the cortex we move from the primary auditory cortex
149 (pmHG) laterally to PT and mSTG, and from there to pSTG, aSTG, and IFG. The corresponding
150 model depth (layer number) of an electrode is also correlated with various other metrics
151 associated with being downstream in speech processing – relative prediction improvement from
152 ASR over baseline (Spearman's $r = 0.582$, $p \ll 1e - 4$), neural response latency ($r = 0.373$,
153 $p \ll 1e - 4$), neural site's distance from the primary auditory cortex (center of pmHG chosen as
154 reference for primary AC; $r = 0.613$, $p \ll 1e - 4$). These results together suggest a clear
155 association between the biological and artificial processing pathways.

156 The brain results, however, show a difference in encoding between hemispheres, where neural
157 sites corresponding to deeper layers of the ASR model are predominantly found in the left
158 hemisphere. We can see this by a statistical comparison between the distribution of associated
159 model depth of electrodes (one-sided two-sample t-test, $t = 6.31$, $p \ll 1e - 4$, $df = 289$).
160 More interesting is the difference by anatomical region, where as we move further away from
161 the primary AC, lateralization increases – pmHG ($t = -0.183$, $p = 0.57$, $df = 26$), aHG ($t =$

162 $-1.25, p = 0.89, df = 64$), PT ($t = 2.57, p = 0.008, df = 25$), mSTG ($t = 2.27, p = 0.0134,$
163 $df = 59$), pSTG ($t = 2.78, p = 0.0039, df = 46$), aSTG (no speech-responsive sites in the right
164 hemisphere), IFG ($t = 2.98, p = 0.0024, df = 40$). These results are in line with previous
165 studies showing left-lateralization of higher-order regions in linguistic processing (21, 22): lower-
166 level acoustic processing (according to ASR) is similar between hemispheres, while higher-level
167 processing is biased towards the left hemisphere.

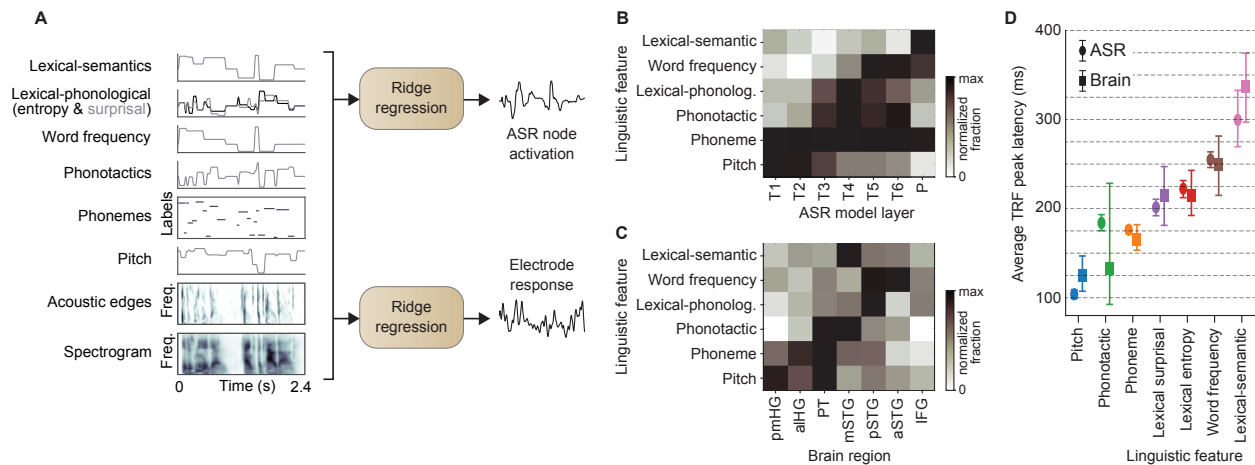
168 **Similar node-level linguistic encoding across the brain and ASR**

169 We used a regression-based method with temporal receptive fields (TRFs) (23) to measure the
170 degree of encoding of different levels of linguistic information in individual nodes of the ASR
171 model and individual neural sites (Figure 2A). In this method, we first predicted the neural activity
172 of a site from a list of acoustic and linguistic features. Next, we determined the contribution of
173 each feature in predicting the response by replacing it with a control variable and interpret the
174 drop in predictive power as the contribution of that feature to the prediction. We repeated the
175 replacement process for each linguistic feature 100 times and measured the t-statistic between
176 actual and control features to determine the significance of the feature.

177 Figure 2B shows the normalized fraction of nodes in each ASR layer that significantly encode each
178 linguistic feature group (fractions divided by sum across feature/row). Figure 2C shows the same,
179 but for normalized fraction of neural sites in each brain region. Figure S1 show the average t-
180 statistic without thresholding, normalized by the sum of each row (feature). We observed that
181 the brain results show broader encoding of higher-level linguistic features in the nonprimary
182 auditory cortex (PT and STG) compared to the primary auditory cortex (HG). This trend is also
183 mirrored in the ASR model. These trends do not depend on the specific choice of the threshold
184 value, as we can see a similar result using average t-values of groups instead of fractions.

185 Figure 2D shows the temporal order of the encoding of these linguistic features in the ASR model
186 and the brain. To measure the encoding latency for each feature, we first select the group of
187 nodes/electrodes that significantly encoded that feature. For each node/electrode, we compute
188 the latency of the peak absolute value of the TRF weight associated with that feature. We average
189 the latencies to obtain a single latency value for each feature. Finally, using bootstrapping, we

190 obtain 95% confidence intervals for the average latency of each feature. In line with prior findings
 191 in the human brain (24), the results show a particular temporal order of emergence for the
 192 different levels of linguistic information in both systems.



193

194 **Figure 2. Similar node-level linguistic encoding across the brain and ASR.**

195 **(A)** Determining node-level linguistic encoding. To measure which linguistic features are

196 encoded in a given electrode site or ASR node, we use a similar approach to ¹. We first predict

197 the neural response r_S to the stimuli S from a set of time-aligned acoustic-linguistic features

198 $\mathcal{F}_S = \{f_{j,S}\}$, where $f_{j,S}$ is the j -th linguistic feature of the stimulus S . Then to find if feature f_j is

199 encoded in the neural response, we replace it with \hat{f}_j and measure the difference between

200 prediction score between the true and control cases. We repeat this process 100 times and

201 compute a statistic of significance of encoding. **(B and C)** Node-level linguistic encoding in ASR

202 layers and brain regions. The values denote fraction of nodes/electrode in a layer/region

203 significantly encoding a given feature. Fractions are normalized by dividing each row by the

204 maximum. **(D)** Temporal order of node-level linguistic encoding in the brain and ASR. Error bars

205 indicate 95% confidence intervals obtained using bias-corrected and accelerated bootstrap

206 (BCa; $N = 1000$) on the mean latency of encoding across nodes/electrodes.

207

208

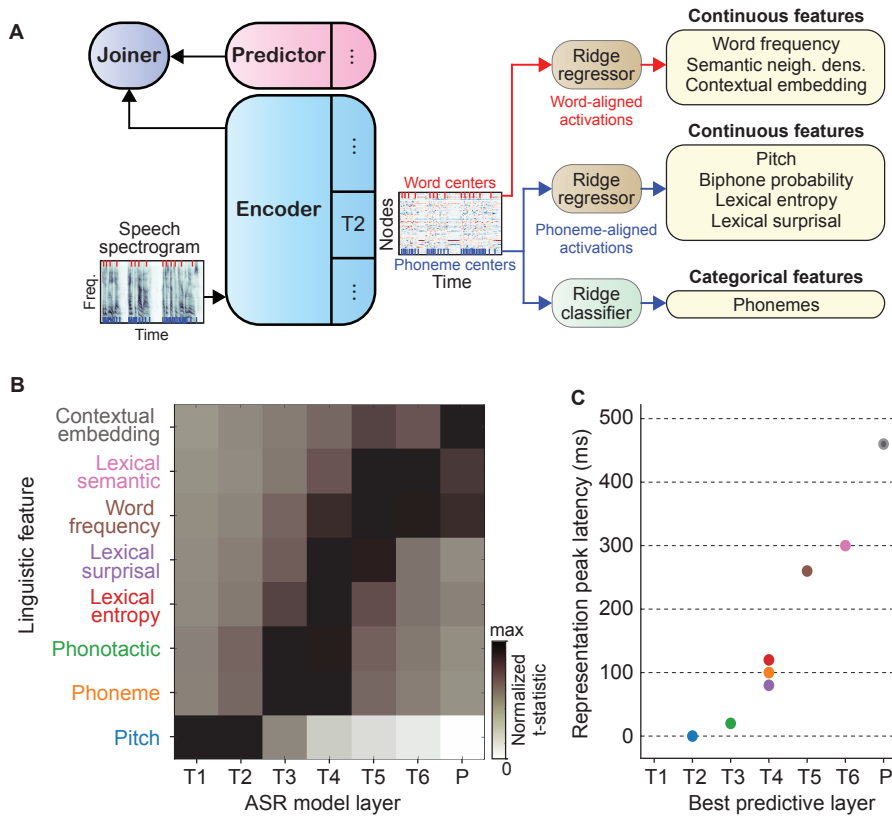
209 **Population-level linguistic encoding in the ASR**

210 While the regression approach enabled us to compare the results between the two systems, it
211 cannot capture population-level representations. This is especially relevant in the case of the ASR
212 model, since we have full access to the entire neural population. To determine the patterns of
213 linguistic encoding in the model, we decoded each linguistic feature from the activations of the
214 model (Figure 1A). Because linguistic features are defined at different resolutions and linguistic
215 units have different lengths, we performed a unit-aligned analysis. For example, to predict
216 linguistic features that are defined at phoneme resolution—pitch, phoneme, biphone probability,
217 lexical entropy and surprisal—we extracted a layer’s activations at time points $\{p_i + \delta\}$, where
218 $\{p_i\}$ were the time points corresponding to phoneme centers and δ was a constant time lag.
219 Similarly, for word resolution features—word frequency, semantic neighborhood density,
220 contextual embedding—we extracted activations at time points $\{w_i + \delta\}$, where $\{w_i\}$ were the
221 time points corresponding to word centers and δ was a constant time lag.

222 For predicting the categorical phoneme labels we fit linear Ridge classifiers, and for the rest of
223 the features, Ridge regressors (24). The best time lag (δ) and regularization parameter (λ) were
224 determined independently for each layer-feature pair by maximizing the 10-fold cross-validated
225 prediction scores.

226 We also decoded each linguistic feature from the 256-dimensional baseline predictor described
227 earlier. Figure 3B shows the prediction scores of predicting each feature from each layer of the
228 model, normalized per feature such that zero corresponds to the baseline prediction score and 1
229 corresponds to the best score across all layers of the model. We can see that as we move deeper
230 into the model, representations of higher-order linguistic information emerge. Based on these
231 scores, we can associate each linguistic feature with a layer of the model and find the best lag for
232 that layer-feature pair. As a result, we can describe the place of a linguistic feature in the speech
233 recognition process by its time lag and layer depth (Figure 3C). Together, these results show a

234 pronounced emergent linguistic representation in the model through time and space and with a
 235 specific order, enabling a direct comparison with the results observed in the brain.



236

237 **Figure 3. Population-level linguistic encoding in the ASR.**

238 **(A)** Decoding linguistic information in the ASR model. To measure how strongly linguistic
 239 features are encoded at the population-level (layer embedding), we decode the features from
 240 the activations of the different layers of the ASR. We use a single-lag Ridge regressor or
 241 classifier. **(B)** Population-level linguistic encoding in different layers of the model. Shown values
 242 are min-max normalized between baseline (spectrogram) prediction score and maximum
 243 prediction score across layers. **(C)** Temporal order of population-level linguistic encoding in the
 244 ASR. For each linguistic feature, we find the layer of the model that best predicts it, then

245 similarly determine the best time lag where the feature was maximally predictable (see
246 Methods).

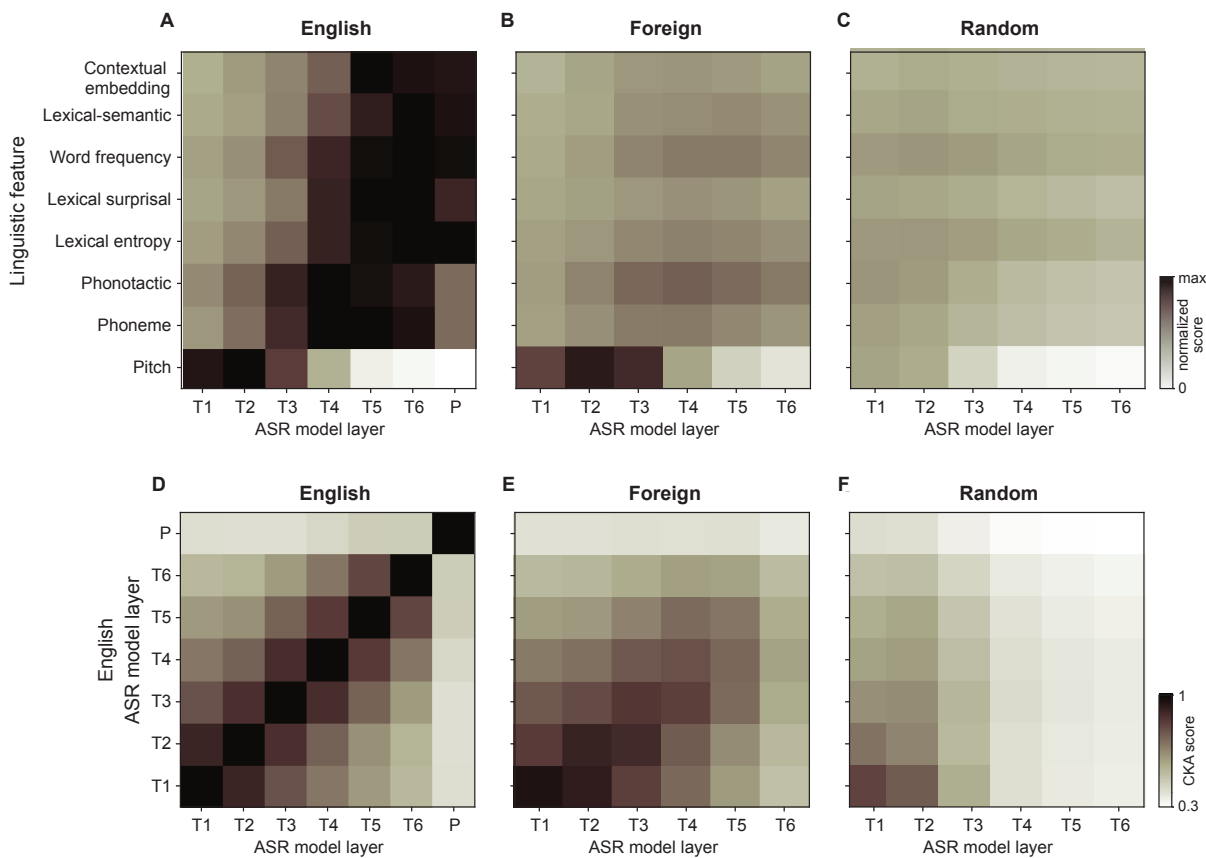
247

248

249 **Effect of model training on linguistic representation**

250 Finally, we explored the conditions in which such a linguistic representation can emerge in the
251 model. We tested two hypotheses: the emergent representation is training-dependent, and the
252 emergent representation is language-dependent. To do this, we decoded the same list of English
253 linguistic features from three different bi-directional RNN-T models: one trained on English, one
254 on French, and a randomly initialized model which underwent no training (Figures 4A–C). We
255 used bi-directional models for this analysis because we do not directly compare the models with
256 the brain, and bi-directional models perform better and are more commonly used in ASR tasks.
257 Additionally, the layers of the bi-directional models have more consistent latency which makes
258 direct comparison between them more accurate, as we do next. The results show a stark
259 difference between the English and random models, such that the random model does not show
260 a strong linguistic representation for any feature compared to the baseline, especially for deeper
261 layers. The French model shows an identical representation of pitch compared to the English
262 model and a relatively strong representation of only phonemes, phonotactics, and word
263 frequency compared to the random model.

264 We can also compare these models more directly by computing the similarities of their
265 representations to the same stimulus using centered kernel alignment (CKA) (25). We compared
266 the representations of each layer of the English model to all layers of all three models (Figures
267 4D–F). We observed that the English and French models share similar representations of their
268 first 3 layers to a large degree, but steadily decrease in similarity over layers overall. In contrast,
269 the random model has weak similarity to the English model, even from early layers. Put together,
270 these results show that the emergent linguistic representation observed in the RNN-T is both
271 training-dependent and language-dependent.



272

273 **Figure 4. Effect of model training on linguistic representation.**

274 **(A–C)** Linguistic encoding in control ASR models. We perform the same analysis as in Figure 3
275 on three bi-directional RNN-Transducers—a model trained to recognize English speech, a model
276 trained to recognize French (“foreign”), and a randomly initialized model (“random”). For the
277 foreign and random models we discard the prediction branch (P) activations, as they fail at
278 transcribing speech so their top-down branch cannot be aligned to the bottom-up branch.
279 Values are min-max normalized between baselined (spectrogram) predictions score and
280 maximum prediction score across all layers of all models. **(D–F)** Pairwise representation
281 similarity analysis between ASR model layers. Each of the three panels shows the centered

282 kernel alignment (CKA) similarity between the activations of each of its layers and each of the
283 layers of the English model.

284

285

286 **Discussion**

287 This study presents a detailed comparison of the linguistic representations in the human brain
288 and an ASR model, specifically focusing on the hierarchical encoding of speech. By using a
289 biologically plausible RNN-Transducer model, we demonstrate similarities and differences in how
290 linguistic features are processed across systems. By selecting a model architecture that mirrors
291 the causal and incremental nature of human speech processing, and incorporating a biologically
292 plausible RNN architecture which processes the input at the spectrogram-level rather than the
293 word-level, we provide a fair ground for comparing how each system encodes various levels of
294 linguistic information, separating this study from prior works using transformer-based large
295 language models or non-causal architectures (6–11). Our findings show a striking alignment
296 between specific brain regions and ASR model layers. Both systems exhibit a hierarchical
297 progression in the encoding of linguistic information, moving from low-level acoustic analysis to
298 higher-level semantic understanding. The correlations between layers and neural regions
299 reconstructs the dorsal and ventral streams of speech processing (26–28), with increasing
300 alignment with deeper layers. Specifically, we find increasing correlation with deeper layers of
301 the model along the dorsal pathway from mSTG to pSTG and finally to IFG, as well as the ventral
302 pathway from mSTG to aSTG and IFG. The fact that IFG most strongly matches the later layers
303 and particularly the prediction branch of the ASR model, which receives feedback from prior
304 unified representations, suggesting its role in unifying lexical information from previous stages
305 into a single high-level representation, as has been suggested by prior research (29). Our findings
306 of the brain’s hierarchical representations support prior studies of the neural hierarchy of
307 linguistic representations (26, 30–32), and we advance this knowledge by highlighting the specific
308 regions and representational latencies at each stage of this hierarchy which are common across
309 both the human brain and the ASR model. This functional congruence underscores the shared

310 mechanisms of linguistic processing in both the brain and ASR systems, revealing that artificial
311 systems trained to emulate human behavior can mirror aspects of biological speech recognition
312 across the full hierarchy of linguistic feature extraction.

313 Our analysis of ASR models with foreign language training or no training also yielded new insights
314 into how the linguistic feature hierarchy forms. We showed that in the French-trained model,
315 only acoustic features like pitch were strongly represented, with higher-level linguistic
316 representations of the English stimuli being weakly encoded, and features at all levels being
317 weakly encoded in the random model. When comparing the representations of the models to
318 each other, we found that the layer representations of the French model quickly diverged from
319 those of the English model after the first two layers, demonstrating where the effects of training
320 begin to compound. One recent study compared the processing of Chinese and English speech
321 by unsupervised transformer speech representation models and the brain (6). They found that
322 models trained on one language better predicted neural responses recorded in STG from native
323 speakers of that same language, whether English or Chinese, and that this performance
324 difference increased over the model's layers. Our results support this and offer a more granular
325 explanation by analyzing the encoding of higher-level linguistic features in biologically-plausible
326 models and over deeper regions in the brain's language processing pathway, showing that
327 models trained on different languages (e.g. French) or untrained models do not exhibit the same
328 level of linguistic representation for English speech as a native English model. These findings
329 suggest that high-level linguistic representations are not simply a byproduct of deeper layers
330 processing longer segments of speech (33) but are instead tightly linked to the model's training.
331 A closer look at the training process, which is akin to language acquisition in humans, could be
332 utilized to study the development of spoken language processing in infants or second language
333 learners.

334 Despite these similarities, the RNN-T ASR model lacks the brain's top-down feedback connections
335 and flexible, multi-directional communication between regions (34, 35). While the ASR model
336 processes speech in a non-lateralized fashion, our results show that the brain shows a strong left-
337 hemisphere bias for higher-level linguistic tasks, in accordance with prior studies (21, 22). This
338 lateralization highlights a key divergence between the two systems, which the current ASR

339 architecture cannot fully account for. These architectural differences suggest that while ASR
340 models can emulate aspects of human speech processing, more biologically accurate models are
341 needed to capture the full complexity of the brain's hierarchical and interactive mechanisms.
342 Future models incorporating binaural input, top-down feedback, and multitask objectives could
343 provide deeper insights into the brain's speech processing pathways. Additionally, ASR models
344 are trained with only a speech recognition loss, forcing them to develop speaker-independent
345 features which likely remove the speaker-specific information such as pitch, matching our finding
346 that pitch is only strongly encoded in the early layers. This may influence the hierarchy of
347 representations formed over layers and represents a point of divergence between ASR models
348 and the brain (36, 37). Nonetheless, the insight that a high-performing ASR system represents
349 linguistic features, from low-level acoustic to high-level contextual features, in a similar
350 hierarchical manner to the brain suggests that similar recent parallels shown between language
351 models and the brain (7–11) can arise in a system that is performing the same overall task as a
352 human listener, from input sound waveform to output words.

353 This study demonstrates a detailed correspondence between linguistic representations in the
354 human brain and an RNN-Transducer model. The RNN-T's causal and incremental processing
355 mirrors key aspects of human speech processing, making it a reasonable proxy for investigating
356 brain-like hierarchical encoding. While our findings are specific to this model, they likely reflect
357 broader principles of linguistic processing shared across systems. The hierarchical progression of
358 encoding observed here, from low-level acoustic to high-level semantic features, is consistent
359 with patterns identified in prior studies using other architectures (38–40), suggesting robustness
360 in the hierarchical nature of linguistic representation even if layer-to-region mappings vary.
361 Crucially, our findings are grounded in empirical data from the human brain, providing a
362 benchmark for comparison and aligning with known principles of neural processing, such as the
363 ventral and dorsal streams of speech (27–29) and left-hemisphere lateralization (21, 22). Future
364 work should explore whether these similarities extend to other biologically plausible ASR
365 architectures to disentangle architecture-specific phenomena from universal principles. While
366 model-specific details may vary, this study highlights that biologically inspired architectures can

367 reveal shared computational strategies between the brain and artificial systems, offering insights
368 into the constraints and mechanisms underlying human and machine language processing.

369 In conclusion, this study bridges neuroscience and artificial intelligence by revealing how both
370 systems process speech in hierarchical layers. Our findings not only advance our understanding
371 of the brain's mechanisms for linguistic encoding but also offer valuable guidance for the
372 development of more sophisticated and biologically informed ASR models. Moving forward,
373 integrating more brain-like features into artificial systems may unlock new possibilities for
374 understanding and enhancing human and machine communication.

375

376 **Methods**

377 **Data**

378 **Participants, neural data, task, and stimuli**

379 Fifteen human patients (7 female, mean age: 36, range: 19-58) with drug-resistant epilepsy were
380 studied. All patients were implanted with intracranial electroencephalography (iEEG) electrodes
381 for epileptogenic foci localization. Twelve of the patients has stereoelectroencephalographic
382 (sEEG) depth electrodes, while the other three had both depth electrodes and subdural grids
383 and/or strips. All recordings were inspected by an epileptologist to ensure they were free of
384 interictal spikes. The patients provided written, informed consent to participate in the research
385 study prior to implantation, and the protocol was approved by the institutional review board at
386 the Feinstein Institutes for Medical Research.

387 The subjects listened to approximately 30 minutes of stories spoken by voice actors. Occasional
388 pauses in the story were added, resulting in 53 trial segments, and the subjects were asked to
389 repeat the most recent sentence before the pause to ensure they were paying attention. All
390 subjects were able to repeat the sentences without issue. iEEG signals were acquired at 3 kHz
391 sampling rate, and the envelope of the high-gamma band (70-150 Hz) was extracted with the
392 Hilbert transform (41). This signal was then z-scored and resampled to 100 Hz.

393 **Electrode selection**

394 Electrodes were projected onto the nearest cortical surface. We selected all electrode sites
395 within the auditory cortex (AC; $N = 335$) and the inferior frontal gyrus (IFG; $N = 191$). Figure 1A
396 shows the general location of IFG and the subregions of the auditory cortex on the FreeSurfer
397 average brain (15).

398 To determine whether an electrode site was speech-responsive, we first selected 53 pre-trial
399 silence segments ($[-1\text{ s}, 0\text{ s}]$ relative to segment onset) and using a 200 ms wide non-overlapping
400 moving average window reduced it into 53×5 data points. We then selected 53 post-onset
401 speech segments ($[0.4\text{ s}, 1.2\text{ s}]$ relative to segment onset) and using a similar moving average
402 window reduced it into 53×4 data points. We performed a two-sample t-test between the two

403 distributions for each electrode to obtain a p-value, and then performed a Benjamini-Hochberg
404 (false discovery rate) correction with an alpha of 0.05 to determine the speech-responsive
405 electrodes. This method was observed to be more robust compared to some alternatives when
406 viewing across the entire regions of interest which included electrodes inside and outside the
407 auditory cortex with different response latencies.

408 We limited our analysis to sites that were determined to be speech-responsive ($N = 291/526$).
409 We labeled the neural sites in both hemispheres based on their anatomical region to enable
410 population tests: posteromedial Heschl's gyrus (pmHG; $N_L = 12/15$, $N_R = 16/19$), anterolateral
411 Heschl's gyrus (alHG; $N_L = 34/36$, $N_R = 32/34$), planum temporale (PT; $N_L = 9/12$, $N_R =$
412 $18/20$), middle superior temporal gyrus (mSTG; $N_L = 44/53$, $N_R = 17/18$), posterior superior
413 temporal gyrus (pSTG; $N_L = 27/32$, $N_R = 21/23$), anterior superior temporal gyrus (aSTG; $N_L =$
414 $19/60$, $N_R = 0/13$), and inferior frontal gyrus (IFG; $N_L = 32/119$, $N_R = 10/72$). The electrode
415 locations and their responsiveness are plotted in Figure 1 on the average FreeSurfer brain, where
416 the color indicates whether an electrode is speech-responsive.

417

418 **RNN-Transducer**

419 To model the speech processing mechanism, we use an RNN-Transducer (RNN-T) (12), a
420 recurrent neural network trained to perform automatic speech recognition, i.e., predicting
421 graphemes from the speech spectrogram. The model consists of two branches—transcription
422 and prediction (Figure 1C). The transcription branch, consisting of six uni-directional LSTM layers
423 of 640 nodes each (T1–T6), acts as a bottom-up encoder. It transforms the input speech
424 spectrogram into a representation used to predict the output labels. The prediction branch,
425 consisting of a single uni-directional LSTM layer of 1024 nodes (P), acts as an internal language
426 model. It stores an internal representation of the prior predictions of the model that, together
427 with the encoder's output, are used to make the next prediction. A joint network merges the final
428 transcription branch layer with the prediction branch's embedding and predicts text characters
429 to generate the speech transcription. Bi-directional versions of the model were used for the
430 analysis of training language dependencies since uni-directional models trained on French were

431 not available. The English network outputs are projected to 42 logits, corresponding to 41
432 characters plus BLANK. Similarly, for French, we use 72 output units. The model is trained with
433 the RNN-T loss on the Switchboard and Fisher datasets, which, put together, consist of about
434 2,000 hours of English phone conversations (see (42) for training details). An additional 4-fold
435 data augmentation was applied to the input spectrograms to allow better generalization of the
436 model, bringing the duration of the unique training data to 10,000 hours.

437 We computed layer activations for all layers of the model in response to the same 30-minute
438 stimulus set that the human participants listened to.

439

440 **Time-aligning RNN-Transducer activations to neural data**

441 The processed neural data is sampled at 100 Hz and has shape $T \times C$, where T is the number of
442 time steps and C is the number of electrode channels. The activations of the layers in the
443 transcription branch of the RNN-T models are sampled at 50 Hz and have the shape $\left\lfloor \frac{T}{2} \right\rfloor \times H_T$,
444 where H_T is the number of nodes in each hidden transcription layer. These activation matrices
445 can be trivially aligned with the neural data sampled through a resample with the integer factor
446 2. The activations of the prediction branch, however, have shape $U \times H_p$, where U is the number
447 of output graphemes and H_p is the number of nodes in the hidden prediction layer. These
448 activation matrices cannot be directly aligned with the neural data since they are based on
449 graphemes instead of time. To make this alignment possible, we need to find the most likely
450 “time warping” between output grapheme indices ($1 \leq u \leq U$) and time ($1 \leq t \leq \left\lfloor \frac{T}{2} \right\rfloor$). We run
451 the forward-backward algorithm defined in section 2.4 of Graves (12) on the grapheme output
452 probability lattice produced by the model to obtain this “time warping” between the two
453 sequences. We then use the $\alpha: t \rightarrow u$ alignment to stretch the activation matrix Z_p with shape
454 $U \times H_p$ into a matrix Z'_p with shape $\left\lfloor \frac{T}{2} \right\rfloor \times H_p$ where $Z'_p[t, :] = Z_p[\alpha(t), :]$. Since Z'_p has the same
455 shape as transcription layer activations, we can trivially align it to the neural data.

456

457 **Linguistic features of speech**

458 **Pitch**

459 As a measure of pitch, we computed the pitch contour (F0) of the speech signal using the
460 PyWORLD python package, which is a python wrapper for the WORLD vocoder (43). We then
461 averaged the value across the duration of each phoneme to obtain a phoneme-average pitch.

462 **Phonemes**

463 For phonetic features, we used the categorical (one-hot encoded) representation of ARPAbet
464 phonemes. We chose this because it allows classification in the layer-level analysis.

465 **Phonotactic features**

466 Phonotactics represent phoneme transition probabilities, so we used the biphone probability P_{ab}
467 for phoneme bigram ab :

$$468 \quad P_{ab} = \frac{\text{freq}_{ab}}{\sum_{p \in S} \text{freq}_{ap}}$$

469
470 where S is the set of all ARPAbet phonemes. To compute the frequencies of each biphone, we
471 used the CMU dictionary (<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>) to convert words to
472 phoneme sequences and counted the occurrence of each biphone using the SUBTLEX-US corpus,
473 which is an English word frequency dataset calculated from movie subtitles (44). Since biphone
474 frequencies were calculated from a word frequency dataset and without access to word
475 transition probability information, we counted the first phoneme transition of words separately
476 from non-first phonemes. For example, the biphones for the phrase “red hat” are the following:
477 $/\#r/$, $/re/$, $/ed/$, $/\#h/$ (not $/dh/$), $/hæ/$, and $/æt/$. The frequency of a phoneme bigram represents
478 the degree of exposure of an average native listener to the bigram and measures its probability
479 in natural speech. We purposefully chose a non-position-specific measure of phonotactics (as
480 opposed to the more common approach (45)) to maximally dissociate this effect from lexical
481 processes. This level represents the expectation and surprisal of the listener when hearing a new

482 phoneme, based on the immediate past. This prelexical phonotactics feature could indicate
483 predictive coding mechanisms that operate on the phonemic level (46–49).

484 **Lexical-phonological features**

485 To measure the lexical-phonological effect, we used lexical entropy and surprisal. These values
486 were calculated for each phoneme within a word from the previous phonemes in that word. The
487 surprisal caused by phoneme φ_i , $S(i)$, in word $w=\varphi_1\dots\varphi_k$ indicates the improbability of hearing
488 phoneme φ_i based on the previous $i-1$ phonemes that came before it in the word and is
489 calculated as follows:

$$490 \quad S(i) = -\log_2 \frac{\text{freq}(\text{cohort}_i)}{\text{freq}(\text{cohort}_{i-1})}$$

491 where $\text{freq}(\text{cohort}_i)$ is the summed frequency of all words that start with the phoneme
492 sequence $\varphi_1\dots\varphi_i$. On the other hand, the lexical entropy, $E(i)$, for phoneme φ_i , is the entropy
493 within all words that start with the phoneme sequence $\varphi_1\dots\varphi_i$ (the cohort) (50):

$$494 \quad E(i) = - \sum_{\text{word} \in \{\text{cohort}_i\}} p(\text{word}) \log_2 p(\text{word})$$

495 where $p(\text{word})$ indicates the relative frequency of the word within the cohort. These two
496 parameters together encode the incremental lexical competition among all phonologically
497 consistent candidates as a word is being heard, weighted by their frequency. To compute lexical
498 surprisal for the word-initial phoneme, we assumed a transition from the entire lexicon, i.e., how
499 surprising it is to hear a word starting with phoneme p given the frequencies of all the words in
500 the lexicon.

501 **Lexical-semantic features**

502 To study the encoding of semantic information, we represented each word with its semantic
503 neighborhood density (SND) obtained from the English Lexicon Project (51, 52), which refers to
504 the relative distance between a word and its closest neighbors based on a global co-occurrence
505 model (51, 52). The neighborhood density can encode the degree of activation of semantically
506 related words in the lexicon upon hearing the target word.

507

508 **Contextual embedding**

509 We used the embedding obtained from the last hidden layer of a pre-trained GPT-2 XL model
510 obtained from Hugging Face (53). This 1.5B parameter version of GPT-2, a transformer-based
511 language model was pretrained using a causal language modeling (CLM) objective on English
512 language data. CLM is a training goal where the model predicts the next token in a sequence
513 given its preceding tokens, ensuring that the prediction for a position can only depend on known
514 outputs at previous positions.

515 To associate contextual embedding to words in the analysis data, we concatenated the
516 transcripts for all 53 trials, in the same order the participants heard them. Then, for each word,
517 we gave the model the last token of the target word and the 511 tokens preceding it (if available)
518 to obtain a contextual embedding for that word. The total 512 tokens are half the maximum
519 context window size of GPT-2 XL and correspond to about 3 minutes of context in the experiment
520 data.

521

522 **Predicting neural responses from ASR activations**

523 We predict the neural activity recorded at each electrode from each of the layers of the ASR
524 model (Figure 1). We map ASR layers instead of ASR nodes to electrodes, because the high-
525 gamma envelope of an iEEG electrode represents a readout of the activity of neighboring neural
526 populations consisting of thousands of neurons rather than individual neurons.

527 To predict a neural response from a layer of the model, we fit a linear single-lag Ridge regressor
528 that predicts the electrode activity from the layer activations in response to the same stimuli.
529 Discounting a constant lag term δ , the output at time t is predicted only from the input at time
530 $t + \delta$. In other words, $y = A_i^{+\delta} \beta + \varepsilon$, where y is the $T \times 1$ neural response vector, $A_i^{+\delta}$ is the
531 $T \times H$ activation matrix for layer i of the ASR model with H nodes shifted by δ time steps across
532 its first dimension, β is the $H \times 1$ regression model that maps $A_i^{+\delta}$ to y , and ε is the regression
533 error.

534 We use the Ridge regressor from the python scikit-learn package to fit these models. The optimal
535 lag value (δ) and ridge regularization parameter (λ) are chosen independently for each electrode-
536 layer pair by maximizing the prediction score through 10-fold cross-validation.

537

538 **Predicting neural responses and ASR activations from linguistic features**

539 We used a temporal receptive field (TRF) (23) framework to measure the extent to which
540 different linguistic features are encoded in the neural responses recorded from the brain and
541 activations extracted from the ASR nodes (Figure 2). We used a broad spectrum of acoustic and
542 linguistic features as predictors: spectrogram, acoustic edges, pitch (phoneme-average F0),
543 phonemes, phonotactics ($-\log P_{ab}$), word frequency, lexical-phonological (lexical entropy and
544 lexical surprisal), and lexical-semantics (semantic neighborhood density). The enormous
545 dimensionality of the contextual embedding feature prevented its use in this framework.

546 To produce a distribution of shuffled models for statistics, we replaced each linguistic feature
547 with a shuffled version—one at a time—100 times and measured a t-statistics of encoding, using
548 $t = 19$ as the threshold of significant encoding. Special care was taken for each type of feature
549 to ensure that the shuffling did not alter the influence of lower-level features in the hierarchy, as
550 described below.

551 **Pitch features**

552 We grouped all the words within our 30-minute data based on their number of phonemes. Then
553 we shuffled the word-to-pitch sequence association map within each group.

554 **Phonetic features**

555 We took the CMU pronunciation dictionary (<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>),
556 grouped words by their length measured in phonemes, and then shuffled the word-to-phoneme
557 mapping within each group. As a result, each word will have a consistent pronunciation at ever
558 occurrence, but words that share phonemes will have independent pronunciations, e.g. /kæt/
559 and /bæt/ no longer share two of their three phonemes. We constrained the reassociation to
560 words of same length so that we kept the phoneme alignment information intact and because

561 words of the same length are more similar in frequency of occurrence (i.e. shorter words tend to
562 be more frequent). This is a rather strict control, since shuffling pronunciations with other actual
563 English words maintains the proper syllabic structure for English words.

564 **Phonotactic features**

565 To generate controls for phonotactic features, we shuffled the bigram-to-frequency associations
566 (i.e. the look-up table for bigram frequencies), which means that each bigram was associated
567 with the frequency of a randomly chosen bigram from the true distribution. This control scheme
568 maintained consistency across multiple occurrences of the same bigram. To counter the effect of
569 the separation caused by the first vs. non-first phoneme grouping, we perform the above
570 shuffling separately for first phones (ones starting with #) and non-first biphones, so that any first
571 vs. non-first effect will be maintained in the control, and thus discounted.

572 **Word frequency**

573 We grouped words based on their phoneme length and shuffled the word-to-frequency
574 associations within each group.

575 **Lexical-phonological features**

576 We grouped all cohorts based on the length of their shared phoneme sequence and shuffled the
577 cohort-to-frequency associations within each group. We used this constrained shuffling to keep
578 the effect of secondary information, such as the phoneme position in the word and word length,
579 unchanged. This control scheme also satisfies consistency, i.e. if two words share their first k
580 phonemes, the cohort information for their first k positions would be the same because the same
581 cohorts are mapped to the same information.

582 **Lexical-semantic features**

583 The control for the semantic condition was constructed by grouping words based on their
584 phoneme length and shuffling the word-to-SND associations within each group.

585

586 **Predicting linguistic features from ASR activations**

587 We predicted different linguistic features from each of the layers of the ASR model. Because
588 linguistic features are defined at different resolutions and linguistic units have different lengths,
589 we performed a unit-aligned analysis. For example, to predict linguistic features that are defined
590 at phoneme resolution—pitch, phoneme, biphone probability, lexical entropy and surprisal—we
591 extracted a layer’s activations at time points $\{p_i + \delta\}$, where $\{p_i\}$ were the time points
592 corresponding to phoneme centers and δ was a constant time lag. This $N \times H$ activation matrix,
593 where N is the number of phonemes and H is the number of nodes, was multiplied by an $H \times D$
594 linear decoder to predict the D -dimensional linguistic feature. Similarly, for word resolution
595 features—word frequency, semantic neighborhood density, contextual embedding—we
596 extracted activations at time points $\{w_i + \delta\}$, where $\{w_i\}$ were the time times corresponding to
597 word centers and δ was a constant time lag.

598 For predicting the categorical phoneme labels we fit linear Ridge classifiers, and for the rest of
599 the features Ridge regressors. The best time lag (δ) and regularization parameter (λ) was
600 determined independently for each layer-feature pair by maximizing the 10-fold cross-validated
601 prediction scores.

602

603 **Representation similarity of ASR model layers**

604 We use two methods to compare representations of two ASR layers. The first is to compare their
605 linguistic decoding results. The second is to compare the representations directly with the
606 centered kernel alignment (CKA) method, which is used to measure the similarity between two
607 sets of high-dimensional vectors, especially representations learned by neural networks (25). For
608 this comparative analysis, we used three bi-directional LSTM models with a similar architecture
609 to the causal one described earlier: a model trained on the same English data; a “random” model
610 that has the same architecture as the English model, but randomly initialized and untrained; a
611 “foreign” model that was trained on a French ASR task and was not exposed to English. We chose
612 bi-directional models for the between-model analysis because the layers of the bi-directional
613 models have more consistent latency which makes the CKA layer comparison more accurate.

614

615 **Code availability**

616 Code for preprocessing neural data, selecting responsive electrodes, and creating brain
617 plots is available in the naplib-python package (54).

618

619 **Acknowledgements**

620 This study was funded by the National Institute on Deafness and Other Communication
621 Disorders. S.B. was also supported by the National Institute on Deafness and Other
622 Communication Disorders, R01DC019979. The funders had no role in the study design, data
623 collection and analysis, decision to publish and manuscript preparation.

624

625 **Competing interests**

626 The authors declare no competing interests.

627

628 References

629

- 630 1. A. Stolcke, J. Droppo, Comparing Human and Machine Errors in Conversational Speech
631 Transcription. arXiv arXiv:1708.08615 [Preprint] (2017).
632 <https://doi.org/10.48550/arXiv.1708.08615>.
- 633 2. I. DeWitt, J. P. Rauschecker, Phoneme and word recognition in the auditory ventral
634 stream. *Proc. Natl. Acad. Sci.* **109**, E505–E514 (2012).
- 635 3. C. J. Price, The anatomy of language: a review of 100 fMRI studies published in 2009.
636 *Ann. N. Y. Acad. Sci.* **1191**, 62–88 (2010).
- 637 4. D. Poeppel, The neuroanatomic and neurophysiological infrastructure for speech and
638 language. *Curr. Opin. Neurobiol.* **28**, 142–149 (2014).
- 639 5. Y. Belinkov, J. Glass, “Analyzing Hidden Representations in End-to-End Automatic
640 Speech Recognition Systems” in *Advances in Neural Information Processing Systems*
641 (Curran Associates, Inc., 2017);
642 [https://proceedings.neurips.cc/paper_files/paper/2017/hash/b069b3415151fa7217e8](https://proceedings.neurips.cc/paper_files/paper/2017/hash/b069b3415151fa7217e870017374de7c-Abstract.html)
643 [70017374de7c-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2017/hash/b069b3415151fa7217e870017374de7c-Abstract.html)vol. 30.
- 644 6. Y. Li, G. K. Anumanchipalli, A. Mohamed, P. Chen, L. H. Carney, J. Lu, J. Wu, E. F.
645 Chang, Dissecting neural computations in the human auditory pathway using deep
646 neural networks for speech. *Nat. Neurosci.* **26**, 2213–2225 (2023).
- 647 7. C. Caucheteux, A. Gramfort, J.-R. King, Evidence of a predictive coding hierarchy in the
648 human brain listening to speech. *Nat. Hum. Behav.* **7**, 430–441 (2023).
- 649 8. C. Caucheteux, J.-R. King, Brains and algorithms partially converge in natural language
650 processing. *Commun. Biol.* **5**, 134 (2022).
- 651 9. A. Goldstein, Z. Zada, E. Buchnik, M. Schain, A. Price, B. Aubrey, S. A. Nastase, A.
652 Feder, D. Emanuel, A. Cohen, others, Shared computational principles for language
653 processing in humans and deep language models. *Nat. Neurosci.* **25**, 369–380 (2022).
- 654 10. E. A. Hosseini, M. Schrimpf, Y. Zhang, S. Bowman, N. Zaslavsky, E. Fedorenko, Artificial
655 neural network language models align neurally and behaviorally with humans even
656 after a developmentally realistic amount of training. *BioRxiv*, 2022–10 (2022).
- 657 11. G. Mischler, Y. A. Li, S. Bickel, A. D. Mehta, N. Mesgarani, Contextual feature extraction
658 hierarchies converge in large language models and the brain. *Nat. Mach. Intell.*, 1–11
659 (2024).
- 660 12. A. Graves, Sequence Transduction with Recurrent Neural Networks. arXiv
661 arXiv:1211.3711 [Preprint] (2012). <https://doi.org/10.48550/arXiv.1211.3711>.
- 662 13. S. Ray, J. H. Maunsell, Different origins of gamma rhythm and high-gamma activity in
663 macaque visual cortex. *PLoS Biol.* **9**, e1000610 (2011).
- 664 14. M. Steinschneider, Y. I. Fishman, J. C. Arezzo, Spectrotemporal analysis of evoked and
665 induced electroencephalographic responses in primary auditory cortex (A1) of the
666 awake monkey. *Cereb. Cortex* **18**, 610–625 (2008).
- 667 15. B. Fischl, A. Van Der Kouwe, C. Destrieux, E. Halgren, F. Ségonne, D. H. Salat, E. Busa,
668 L. J. Seidman, J. Goldstein, D. Kennedy, others, Automatically parcellating the human
669 cerebral cortex. *Cereb. Cortex* **14**, 11–22 (2004).

- 670 16. Y. He, T. N. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, D. Zhao, D. Rybach, A.
671 Kannan, Y. Wu, R. Pang, Q. Liang, D. Bhatia, Y. Shangguan, B. Li, G. Pundak, K. C. Sim,
672 T. Bagby, S. Chang, K. Rao, A. Gruenstein, “Streaming End-to-end Speech Recognition
673 for Mobile Devices” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics,
674 Speech and Signal Processing (ICASSP)* (2019);
675 <https://ieeexplore.ieee.org/abstract/document/8682336>), pp. 6381–6385.
- 676 17. K. Rao, H. Sak, R. Prabhavalkar, “Exploring architectures, data and units for streaming
677 end-to-end speech recognition with RNN-transducer” in *2017 IEEE Automatic Speech
678 Recognition and Understanding Workshop (ASRU)* (2017);
679 <https://ieeexplore.ieee.org/abstract/document/8268935>), pp. 193–199.
- 680 18. J. Li, R. Zhao, H. Hu, Y. Gong, “Improving RNN Transducer Modeling for End-to-End
681 Speech Recognition” in *2019 IEEE Automatic Speech Recognition and Understanding
682 Workshop (ASRU)* (2019); <https://ieeexplore.ieee.org/abstract/document/9003906>), pp.
683 114–121.
- 684 19. L. E. Shafey, H. Soltau, I. Shafran, Joint Speech Recognition and Speaker Diarization via
685 Sequence Transduction. arXiv arXiv:1907.05337 [Preprint] (2019).
686 <https://doi.org/10.48550/arXiv.1907.05337>.
- 687 20. M. Ghodsi, X. Liu, J. Apfel, R. Cabrera, E. Weinstein, “Rnn-Transducer with Stateless
688 Prediction Network” in *ICASSP 2020 - 2020 IEEE International Conference on
689 Acoustics, Speech and Signal Processing (ICASSP)* (2020);
690 <https://ieeexplore.ieee.org/abstract/document/9054419>), pp. 7049–7053.
- 691 21. C. McGettigan, S. Evans, S. Rosen, Z. K. Agnew, P. Shah, S. K. Scott, An Application of
692 Univariate and Multivariate Approaches in fMRI to Quantifying the Hemispheric
693 Lateralization of Acoustic and Linguistic Processes. *J. Cogn. Neurosci.* **24**, 636–652
694 (2012).
- 695 22. R. J. Zatorre, A. C. Evans, E. Meyer, A. Gjedde, Lateralization of Phonetic and Pitch
696 Discrimination in Speech Processing. *Science* **256**, 846–849 (1992).
- 697 23. M. J. Crosse, G. M. Di Liberto, A. Bednar, E. C. Lalor, The multivariate temporal
698 response function (mTRF) toolbox: a MATLAB toolbox for relating neural signals to
699 continuous stimuli. *Front. Hum. Neurosci.* **10**, 604 (2016).
- 700 24. M. Keshishian, S. Akkol, J. Herrero, S. Bickel, A. D. Mehta, N. Mesgarani, Joint,
701 distributed and hierarchically organized encoding of linguistic features in the human
702 auditory cortex. *Nat. Hum. Behav.* **7**, 740–753 (2023).
- 703 25. S. Kornblith, M. Norouzi, H. Lee, G. Hinton, “Similarity of neural network
704 representations revisited” in *International Conference on Machine Learning* (PMLR,
705 2019), pp. 3519–3529.
- 706 26. G. Hickok, D. Poeppel, The cortical organization of speech processing. *Nat. Rev.
707 Neurosci.* **8**, 393–402 (2007).
- 708 27. G. Hickok, D. Poeppel, Dorsal and ventral streams: a framework for understanding
709 aspects of the functional anatomy of language. *Cognition* **92**, 67–99 (2004).
- 710 28. A. D. Friederici, The cortical language circuit: from auditory perception to sentence
711 comprehension. *Trends Cogn. Sci.* **16**, 262–268 (2012).
- 712 29. P. Hagoort, On Broca, brain, and binding: a new framework. *Trends Cogn. Sci.* **9**, 416–
713 423 (2005).

- 714 30. U. Hasson, E. Yang, I. Vallines, D. J. Heeger, N. Rubin, A hierarchy of temporal receptive
715 windows in human cortex. *J. Neurosci.* **28**, 2539–2550 (2008).
- 716 31. Y. Lerner, C. J. Honey, L. J. Silbert, U. Hasson, Topographic mapping of a hierarchy of
717 temporal receptive windows using a narrated story. *J. Neurosci.* **31**, 2906–2915 (2011).
- 718 32. N. Ding, L. Melloni, A. Yang, Y. Wang, W. Zhang, D. Poeppel, Characterizing neural
719 entrainment to hierarchical linguistic units using electroencephalography (EEG). *Front.*
720 *Hum. Neurosci.* **11**, 481 (2017).
- 721 33. M. Keshishian, S. Norman-Haignere, N. Mesgarani, “Understanding Adaptive,
722 Multiscale Temporal Integration In Deep Speech Recognition Systems” in *Advances in*
723 *Neural Information Processing Systems* (Curran Associates, Inc., 2021;
724 [https://proceedings.neurips.cc/paper/2021/hash/ccce2fab7336b8bc8362d115dec2d](https://proceedings.neurips.cc/paper/2021/hash/ccce2fab7336b8bc8362d115dec2d5a2-Abstract.html)
725 [5a2-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/ccce2fab7336b8bc8362d115dec2d5a2-Abstract.html))vol. 34, pp. 24455–24467.
- 726 34. C. D. Gilbert, M. Sigman, Brain States: Top-Down Influences in Sensory Processing.
727 *Neuron* **54**, 677–696 (2007).
- 728 35. K. Kveraga, A. S. Ghuman, M. Bar, Top-down predictions in the cognitive brain. *Brain*
729 *Cogn.* **65**, 145–168 (2007).
- 730 36. C. Tang, L. S. Hamilton, E. F. Chang, Intonational speech prosody encoding in the
731 human auditory cortex. *Science* **357**, 797–801 (2017).
- 732 37. Y. Li, C. Tang, J. Lu, J. Wu, E. F. Chang, Human cortical encoding of pitch in tonal and
733 non-tonal languages. *Nat. Commun.* **12**, 1161 (2021).
- 734 38. T. Nagamine, N. Mesgarani, “Understanding the Representation and Computation of
735 Multilayer Perceptrons: A Case Study in Speech Recognition” in *Proceedings of the*
736 *34th International Conference on Machine Learning*, D. Precup, Y. W. Teh, Eds. (PMLR,
737 2017; <https://proceedings.mlr.press/v70/nagamine17a.html>)vol. 70 of *Proceedings of*
738 *Machine Learning Research*, pp. 2564–2573.
- 739 39. T. Nagamine, M. L. Seltzer, N. Mesgarani, “Exploring how deep neural networks form
740 phonemic categories.” in *Interspeech* (2015), pp. 1912–1916.
- 741 40. Q. Raymondaud, M. Rouvier, R. Dufour, Probing the Information Encoded in Neural-
742 based Acoustic Models of Automatic Speech Recognition Systems. arXiv
743 arXiv:2402.19443 [Preprint] (2024). <https://doi.org/10.48550/arXiv.2402.19443>.
- 744 41. E. Edwards, M. Soltani, W. Kim, S. S. Dalal, S. S. Nagarajan, M. S. Berger, R. T. Knight,
745 Comparison of time–frequency responses and the event-related potential to auditory
746 speech stimuli in human cortex. *J. Neurophysiol.* **102**, 377–386 (2009).
- 747 42. G. Saon, Z. Tüske, D. Bolanos, B. Kingsbury, “Advancing RNN Transducer Technology
748 for Speech Recognition” in *ICASSP 2021 - 2021 IEEE International Conference on*
749 *Acoustics, Speech and Signal Processing (ICASSP)* (2021;
750 <https://ieeexplore.ieee.org/abstract/document/9414716>), pp. 5654–5658.
- 751 43. M. Masanori, K. Hideki, K. Haruhiro, Fast and reliable f0 estimation method based on
752 the period extraction of vocal fold vibration of singing voice and speech. *J. Audio Eng.*
753 *Soc.* (2009).
- 754 44. M. Brysbaert, B. New, Moving beyond Kučera and Francis: A critical evaluation of
755 current word frequency norms and the introduction of a new and improved word
756 frequency measure for American English. *Behav. Res. Methods* **41**, 977–990 (2009).

- 757 45. M. S. Vitevitch, P. A. Luce, Probabilistic Phonotactics and Neighborhood Activation in
758 Spoken Word Recognition. *J. Mem. Lang.* **40**, 374–408 (1999).
- 759 46. C. Brodbeck, L. E. Hong, J. Z. Simon, Rapid Transformation from Auditory to Linguistic
760 Representations of Continuous Speech. *Curr. Biol.* **28**, 3976–3983.e5 (2018).
- 761 47. S. Ylinen, M. Huuskonen, K. Mikkola, E. Saure, T. Sinkkonen, P. Paavilainen, Predictive
762 coding of phonological rules in auditory cortex: A mismatch negativity study. *Brain*
763 *Lang.* **162**, 72–80 (2016).
- 764 48. K. Friston, The free-energy principle: a rough guide to the brain? *Trends Cogn. Sci.* **13**,
765 293–301 (2009).
- 766 49. P. Gagnepain, R. N. Henson, M. H. Davis, Temporal Predictive Codes for Spoken Words
767 in Auditory Cortex. *Curr. Biol.* **22**, 615–621 (2012).
- 768 50. M. K. Leonard, K. E. Bouchard, C. Tang, E. F. Chang, Dynamic Encoding of Speech
769 Sequence Probability in Human Temporal Cortex. *J. Neurosci.* **35**, 7203–7214 (2015).
- 770 51. D. A. Balota, M. J. Yap, K. A. Hutchison, M. J. Cortese, B. Kessler, B. Loftis, J. H. Neely,
771 D. L. Nelson, G. B. Simpson, R. Treiman, The English Lexicon Project. *Behav. Res.*
772 *Methods* **39**, 445–459 (2007).
- 773 52. C. Shaoul, C. Westbury, Exploring lexical co-occurrence space using HiDEx. *Behav.*
774 *Res. Methods* **42**, 393–413 (2010).
- 775 53. A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, others, Language models
776 are unsupervised multitask learners. *OpenAI Blog* **1**, 9 (2019).
- 777 54. G. Mischler, V. Raghavan, M. Keshishian, N. Mesgarani, naplib-python: Neural acoustic
778 data processing and analysis tools in python. *Softw. Impacts* **17**, 100541 (2023).

779